

Data Quality Asessment of events data on ORBIT

By Hongxing Niu

2019-10-04

Preface

After you have acquired the data, you should do the Data Quality Assessment (DQA) before to move forward to Advanced Analytics.

The structure of this kernel is as follows

- Preface
- Load data
- Data quality assessment using metrics of six dimensions: DQA1-DQA6
- Automated reporting

Purpose and structure of this kernel

- Give a systematic assessment of the quality of the raw data on ORBIT about events and Events link
 - **1.Validity/Semantic Accurary** *Check if data conforms to the syntax (format, type, range) of its definition*
 - **2.Cardinality/uniqueness/duplicates** *Nothing will be recorded more than once based upon how that thing is identified*
 - **3.Completeness**
 - a. **3.1. Overall completeness** *The overall proportion of stored data against the potential of “100% complete”*
 - b. **3.2. Time series completeness** *The proportion of stored data against the potential of “100% complete” in the course of time*
 - **4.Accuracy**
 - a. **4.1. Distribution of raw data** *Distribution analysis entails counting all the records associated with each value and dividing these by the total number of records to see what percentage of the data is associated with any specific value and how the percentages compare to each other*
 - b. **4.2. Outlier detection** *An outlier is a value that lies in the tail of the statistical distribution of a set of data values(usually ± 2 SD)*
 - **5.Consistency** *The absence of difference, when comparing two or more representations of a thing against a definition*
 - **6.Timeliness/availability** *The degree to which data represent reality from the required point in time*
- Performs a data diagnosis and automatically generates a DQA report.

This document introduces **Data Quality Assessment** methods. You will learn how to diagnose the quality of `tbl_df` data that inherits from `data.frame` and `data.frame` with functions provided by `dlookr`.

This package is in synergy with `dplyr`. Particularly in data exploration and data wrangle, it increases the efficiency of the `tidyverse` package group.

Table 1: A knitr kable

Id	UCB_Business_Unit__c	UCB_External_Id__c	Name
a0pG000000AATC1IAP	IMM	E-00062294	2016-09-17-Bonn-Sponsoring-Patienten-
a0p4A00000A0wciQAB	IMM	E-00106476	2018-04-25 Strategisches Beraterboard I
a0p4A00000EJaNIQA1	IMM	E-00122393	2018-11-27-Essen-Dermaxchange
a0p4A00000A0x9WQAR	IMM	E-00106624	2018-03-26München Round Table Prof.
a0p4A00000A0x9vQAB	IMM	E-00106628	2018-05-30 Hamburg Immunology for yo
a0pG000000BQxm3IAD	IMM	E-00075967	2017-10-07 Hamburg Weltrheumatag Pa
a0p4A000009z7rxQAA	IMM	E-00093724	2017-09-21-EssenProjektbesprechung
a0p4A000009zmEQQAY	IMM	E-00099660	2018-02-16-Würzburg-Sponsoring
a0p4A00000CldOHQAZ	IMM	E-00112519	2018-08-30-Halle-Sponsoring-Sommersch
a0pG000000ARWQnIAP	IMM	E-00065933	2016-11-03 Demmin Expertengespräch

Supported data structures

Data diagnosis supports the following data structures.

- data frame : data.frame class.
- data table : tbl_df class.
- **table of DBMS** : table of the DBMS through tbl_dbi.
- **Using dplyr backend for any DBI-compatible database.**

Load libraries

0. Load Data

```
[1] "Dimension of data:"
[1] 1632  24
```

0.Preview data

```
[1] "Preview of 'Events_nPVU_Germany' data frame"
```

1.DQA Dimension1/Rule1: Validity/Semantic Accurary

- Picking the data subset relevant to Epilepsy other than Parkinson’s disease
- Cleansing up:
- Extract a substring (like “2017-09-21” from “2017-09-21T09:30:29.000+0000”) both for Date and CreateDate;+Mis-format of “Call2_vod__r.CreatedDate“; +Wrongly formatted”Call2_vod__r.UCB_Owner_Full_N and “Call2_vod__r.Account_vod__r.UCB_Primary_Grand_Parent__c”

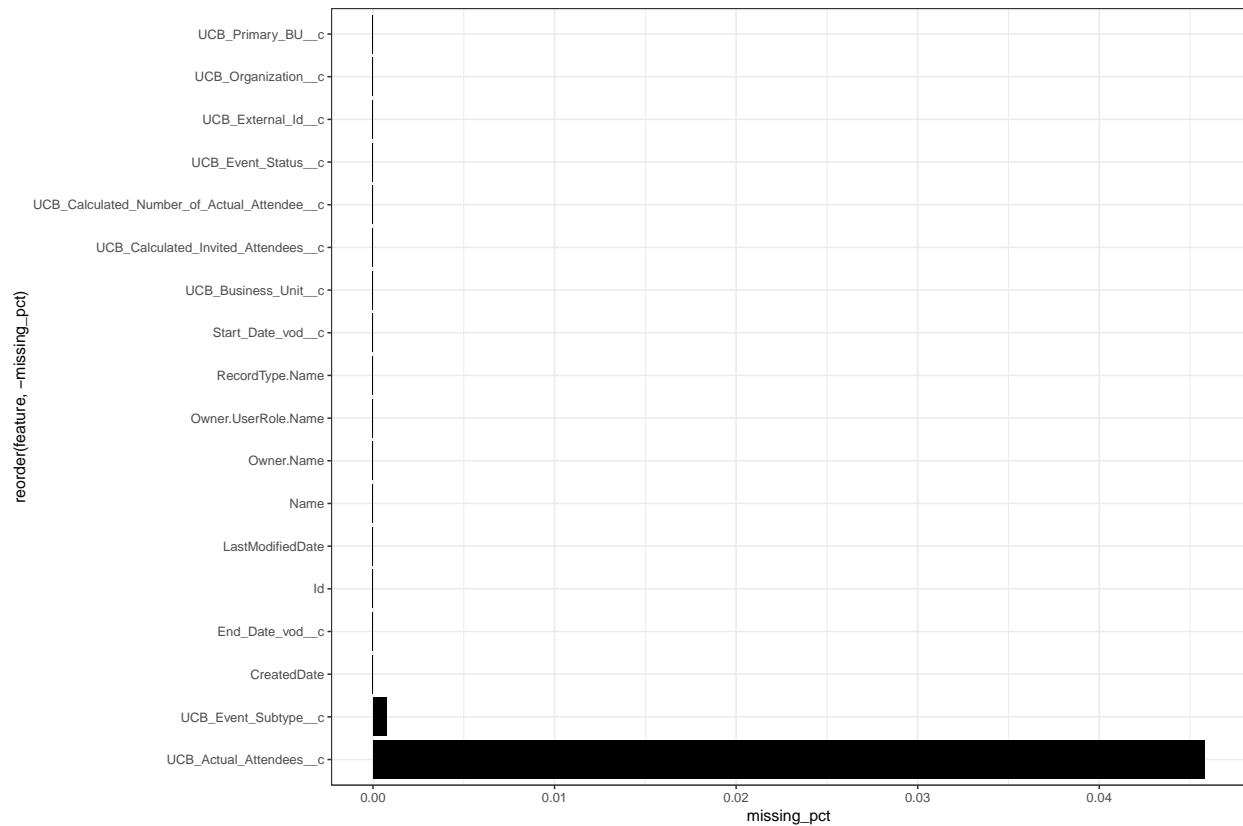
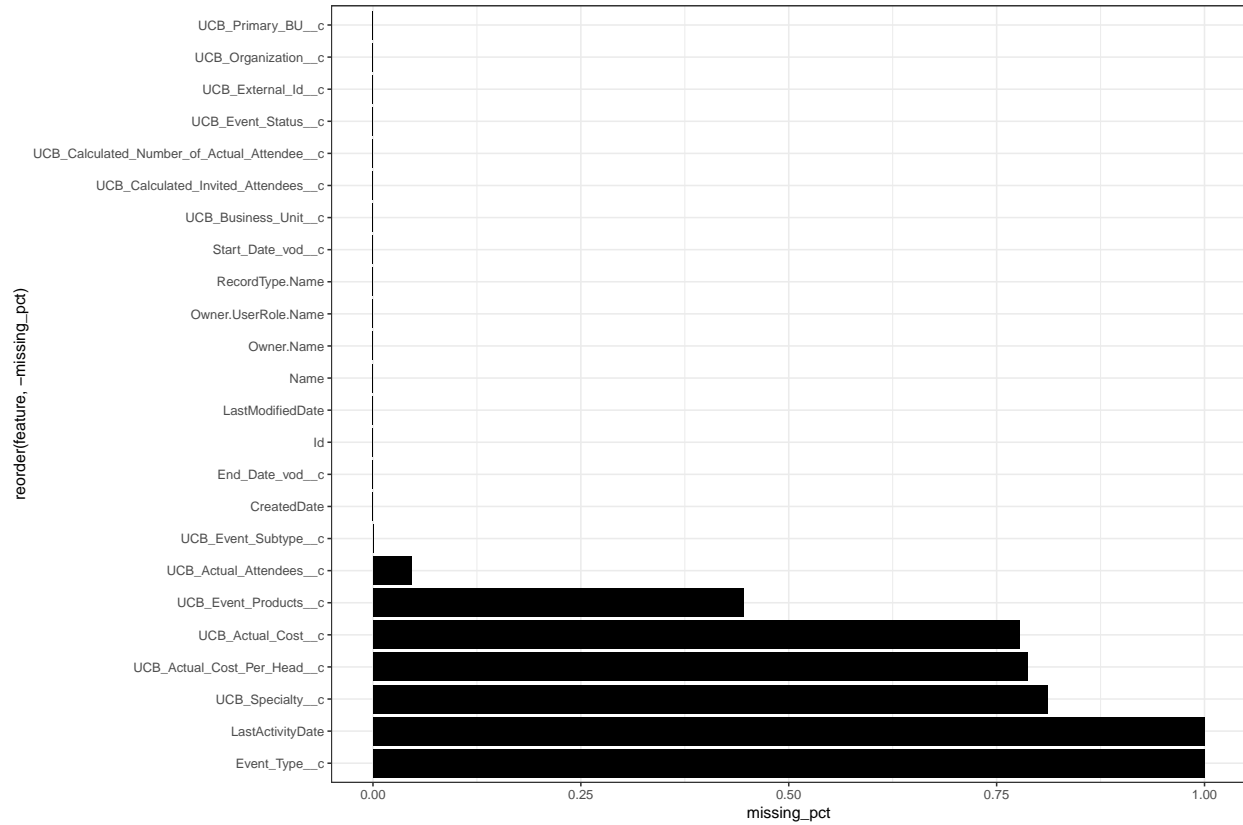
2. DQA Dimension2/Rule2: Cardinality/uniqueness/duplicates

- Based on business understanding of ORBIT data, a hiddent feature “Discussion Priority” differentiates the duplicates. Many duplicates because of “Discussion Priority” behind the data (hinted by “Call2_vod__r.Detailed_Products_vod__c“)

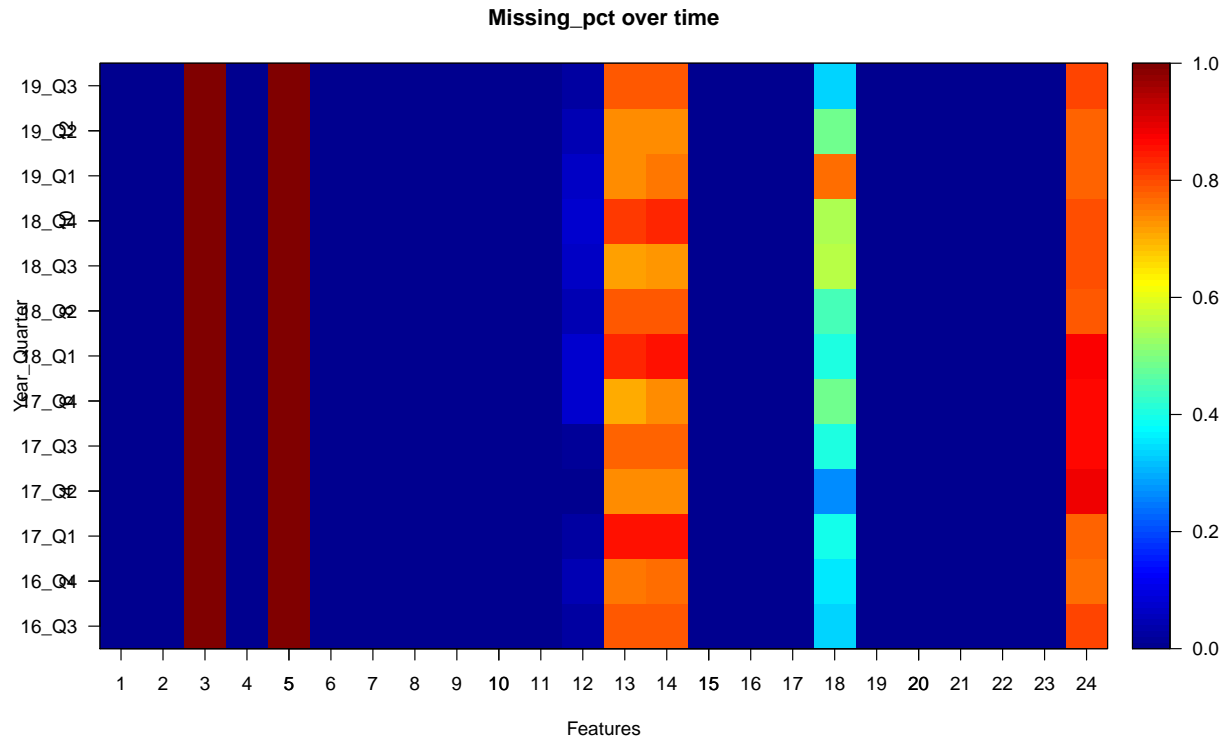
3. DQA Dimension3/Rule3: Handling missing data

We want to get rid of any features that have $\geq 10\%$ of the data missing. We're presuming there's little information in those features, and just adds noise. We plot all the features against the percent of missing values. We keep only those features with $< 10\%$ missing values, and re-adjust the data frame, accordingly.

3.1 find percent missing values for each feature during all period



3.2 Handling historical missing data at at level of quarter, we use 3D plot, with time included to visualize the time series data select the columns to visualize



```
[1] "# Feature"
[1] "CreatedDate"
[2] "End_Date_vod__c"
[3] "Event_Type__c"
[4] "Id"
[5] "LastActivityDate"
[6] "LastModifiedDate"
[7] "Name"
[8] "Owner.Name"
[9] "Owner.UserRole.Name"
[10] "RecordType.Name"
[11] "Start_Date_vod__c"
[12] "UCB_Actual_Attendees__c"
[13] "UCB_Actual_Cost__c"
[14] "UCB_Actual_Cost_Per_Head__c"
[15] "UCB_Business_Unit__c"
[16] "UCB_Calculated_Invited_Attendees__c"
[17] "UCB_Calculated_Number_of_Actual_Attendee__c"
[18] "UCB_Event_Products__c"
[19] "UCB_Event_Status__c"
[20] "UCB_Event_Subtype__c"
[21] "UCB_External_Id__c"
[22] "UCB_Organization__c"
[23] "UCB_Primary_BU__c"
[24] "UCB_Specialty__c"
```

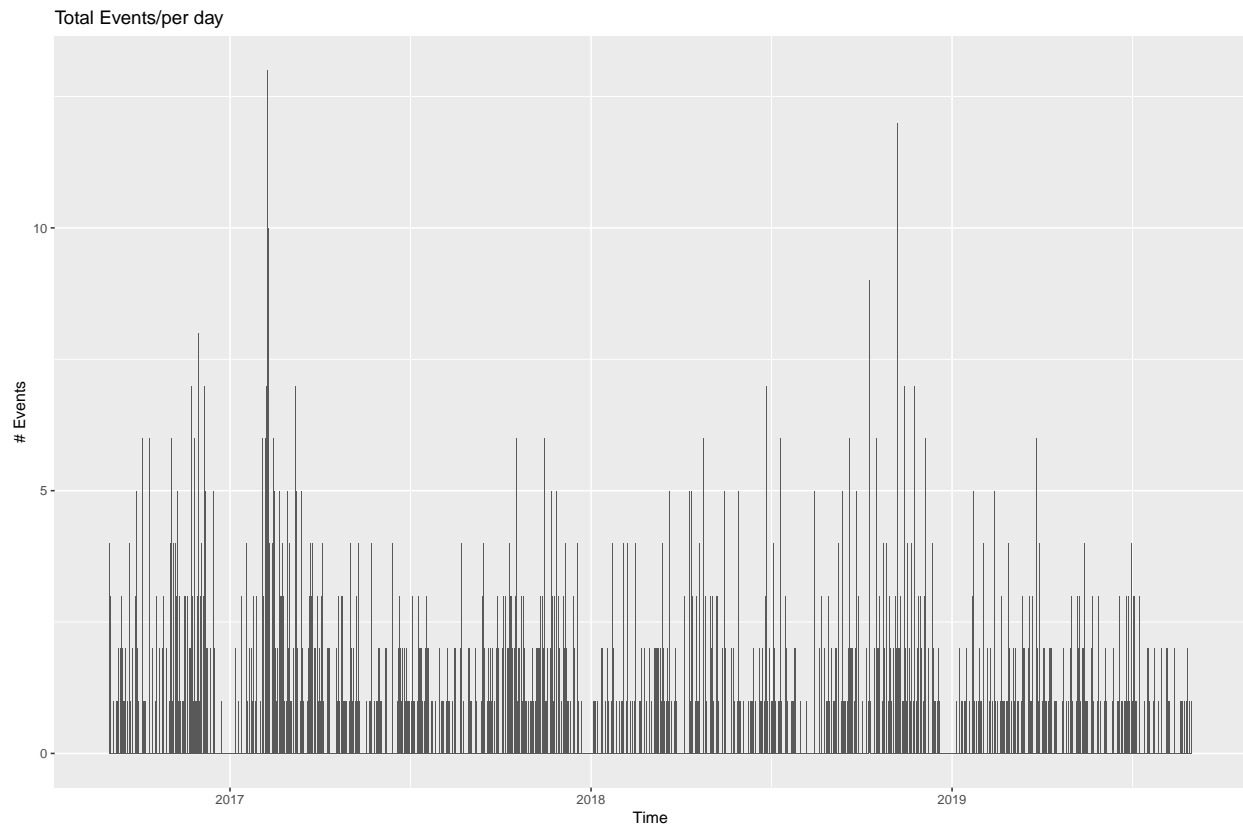
4. DQA Dimension4/Rule4: Distribution and Outlier Detection

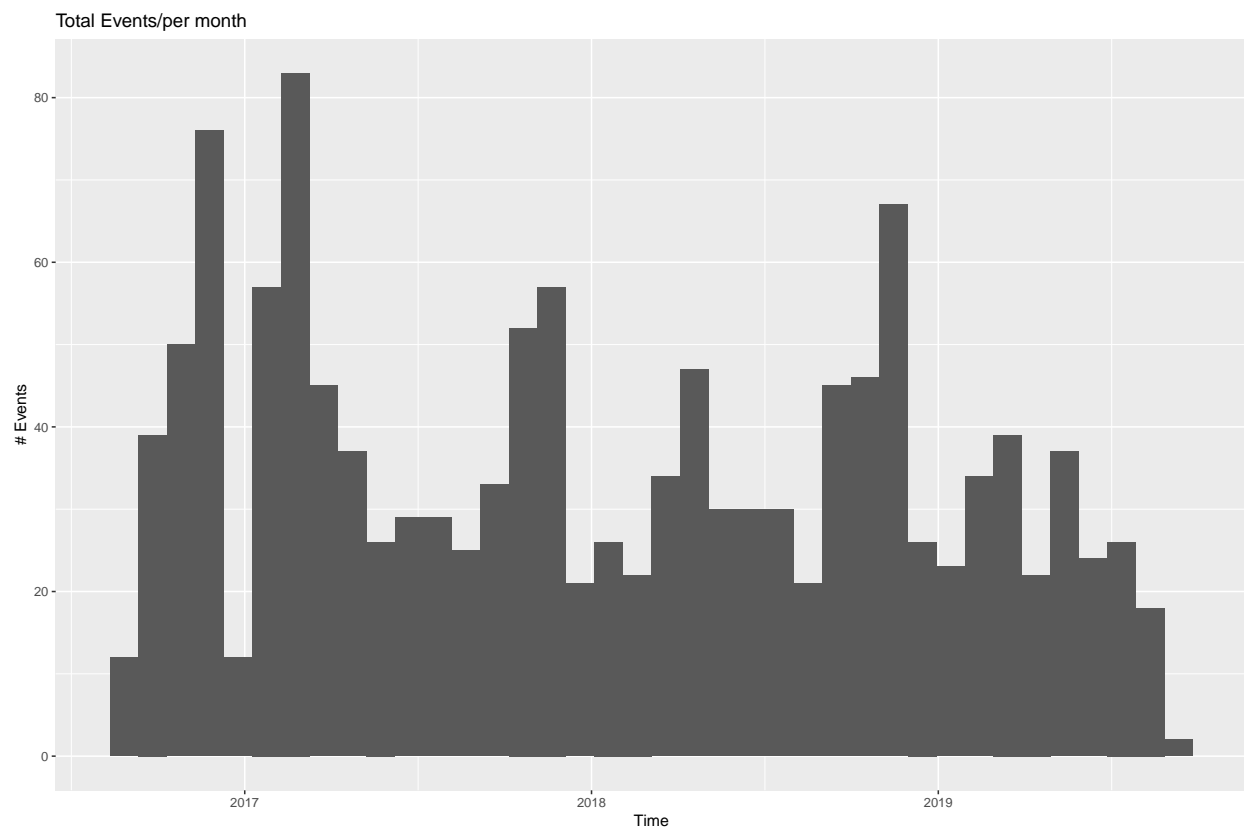
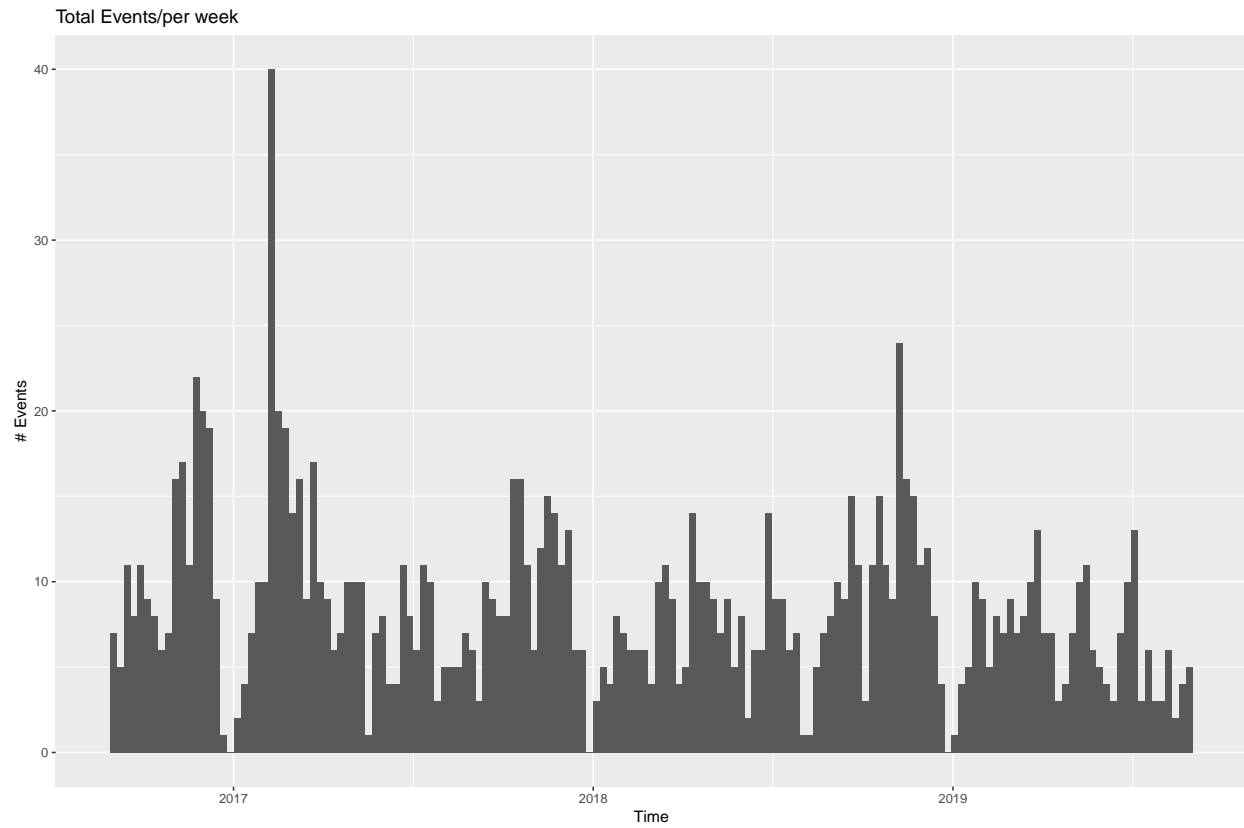
4.1 Distribution of raw data

4.1.1 Visualization and statistics

pdf

2





4.1.2 Normality test on numeric variables using `normality()`

`normality()` performs a normality test on numerical data. Shapiro-Wilk normality test is performed. If the number of observations is larger than 5000, 5000 observations are extracted by random simple sampling and then tested.

The variables of `tbl_df` object returned by `normality()` are as follows.

- `statistic` : Statistics of the Shapiro-Wilk test
- `p_value` : p-value of the Shapiro-Wilk test
- `sample` : Number of sample observations performed Shapiro-Wilk test

We can use `dplyr` to sort non-normal distribution variables by `p_value`:

```
# A tibble: 5 x 4
  vars                statistic p_value sample
<chr>                <dbl>    <dbl>   <dbl>
1 UCB_Actual_Attendees__c 0.108 1.34e-60 1332
2 UCB_Calculated_Invited_Attendees__c 0.671 6.14e-45 1332
3 UCB_Calculated_Number_of_Actual_Attendee__c 0.684 2.51e-44 1332
4 UCB_Actual_Cost_Per_Head__c 0.197 4.97e-33 1332
# ... with 1 more row
```

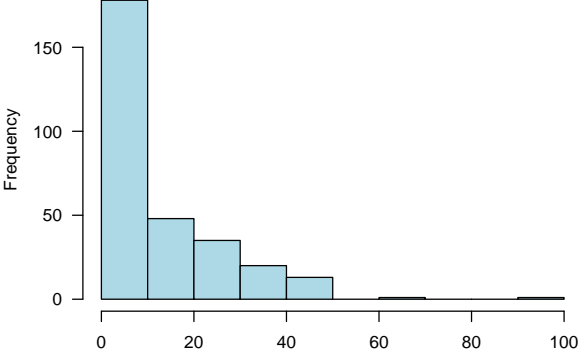
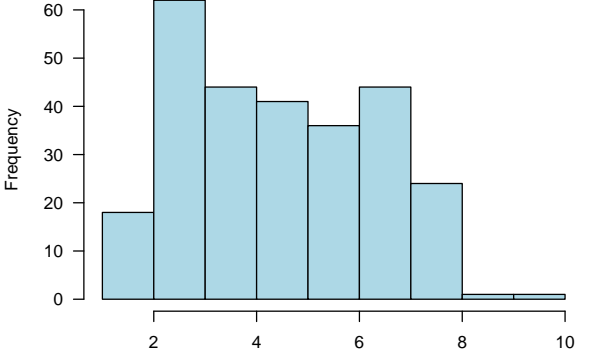
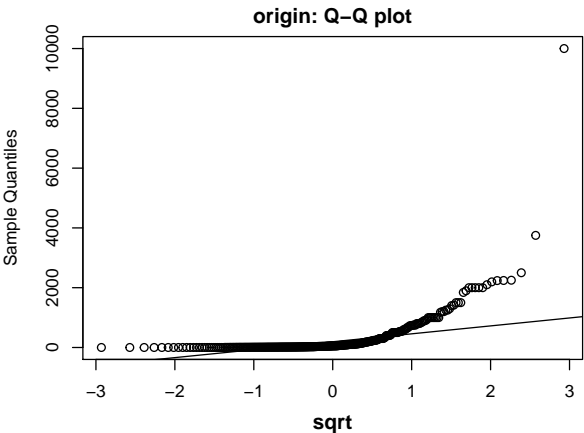
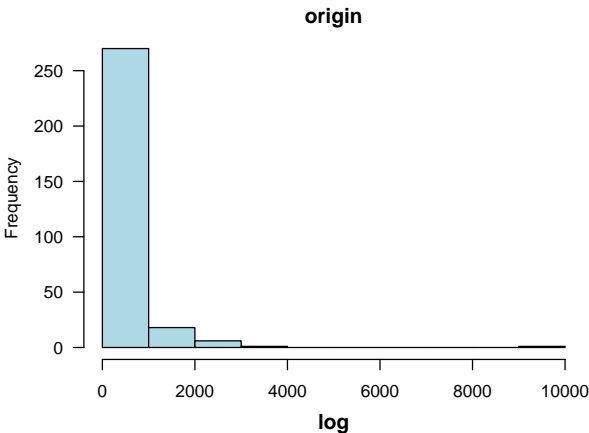
Normalization visualization of numerical variables using `plot_normality()` `plot_normality()` visualizes the normality of numeric data.

The information that `plot_normality()` visualizes is as follows.

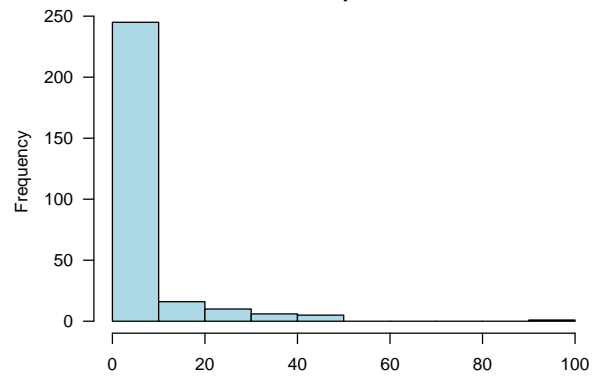
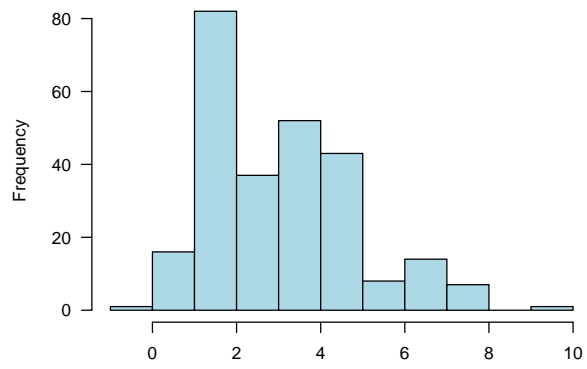
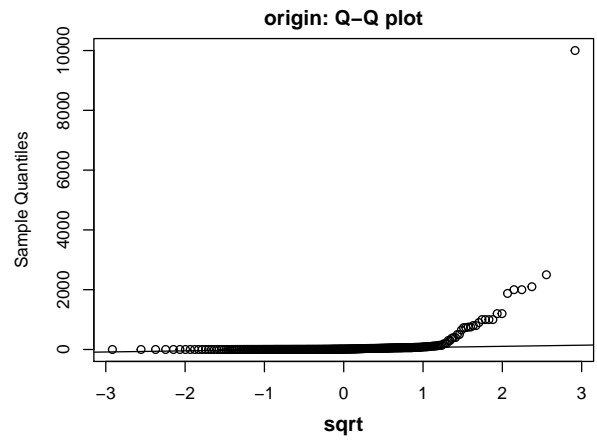
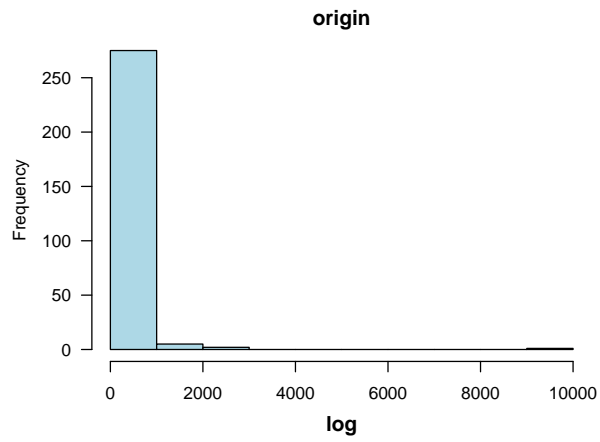
- Histogram of original data
- Q-Q plot of original data
- histogram of log transformed data
- Histogram of square root transformed data

`plot_normality()` can also specify several variables like `normality()` function.

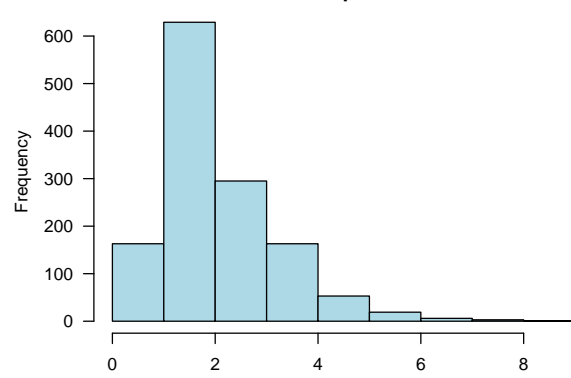
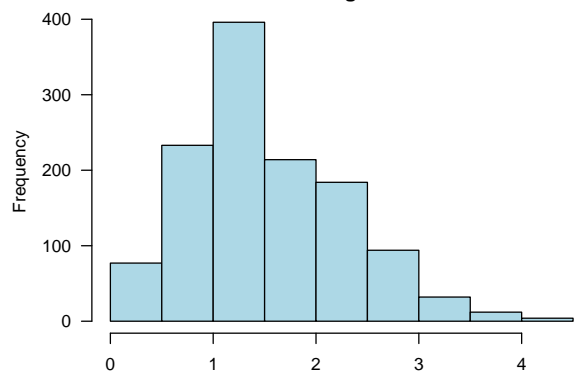
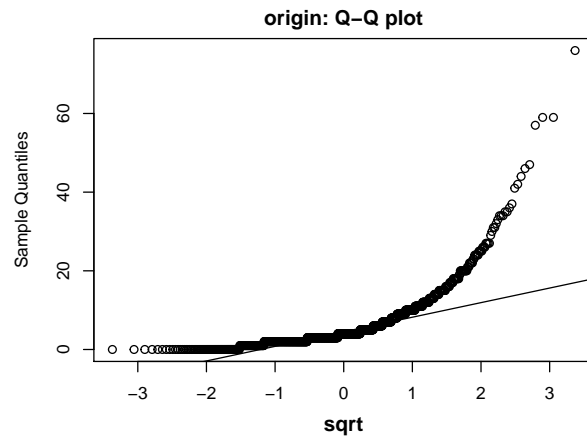
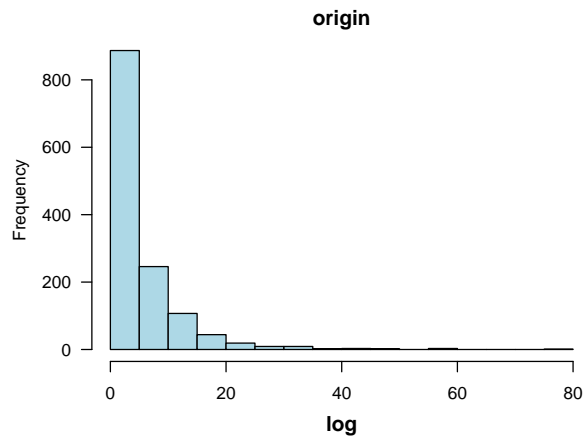
Normality Diagnosis Plot (UCB_Actual_Cost__c)



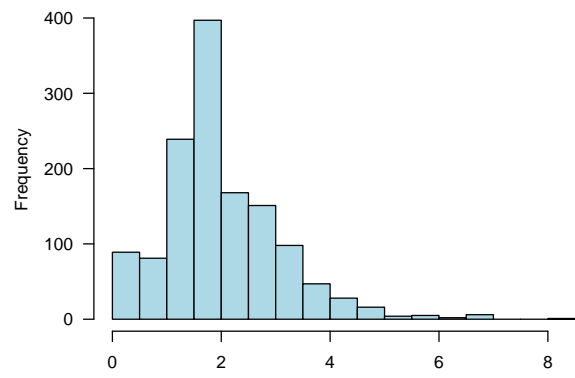
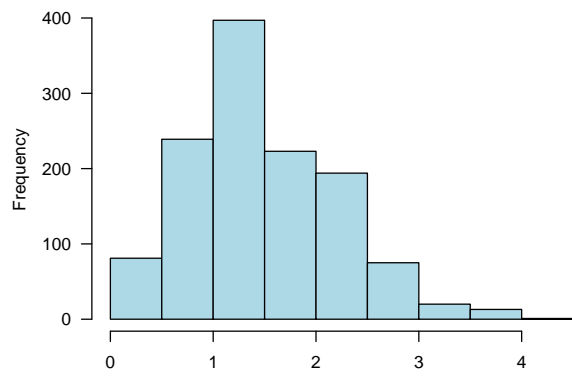
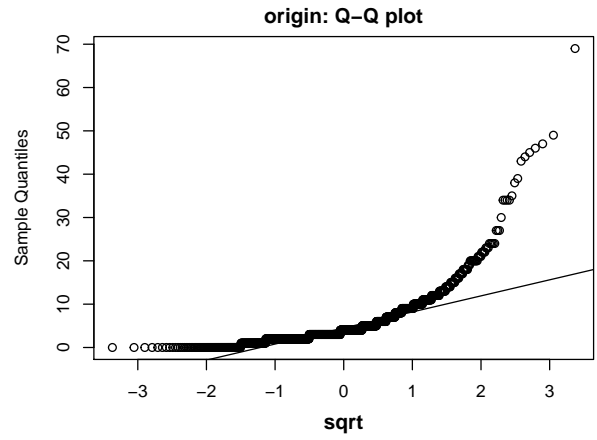
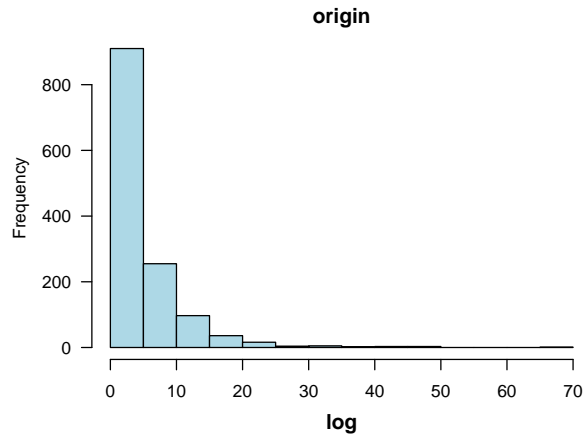
Normality Diagnosis Plot (UCB_Actual_Cost_Per_Head_c)



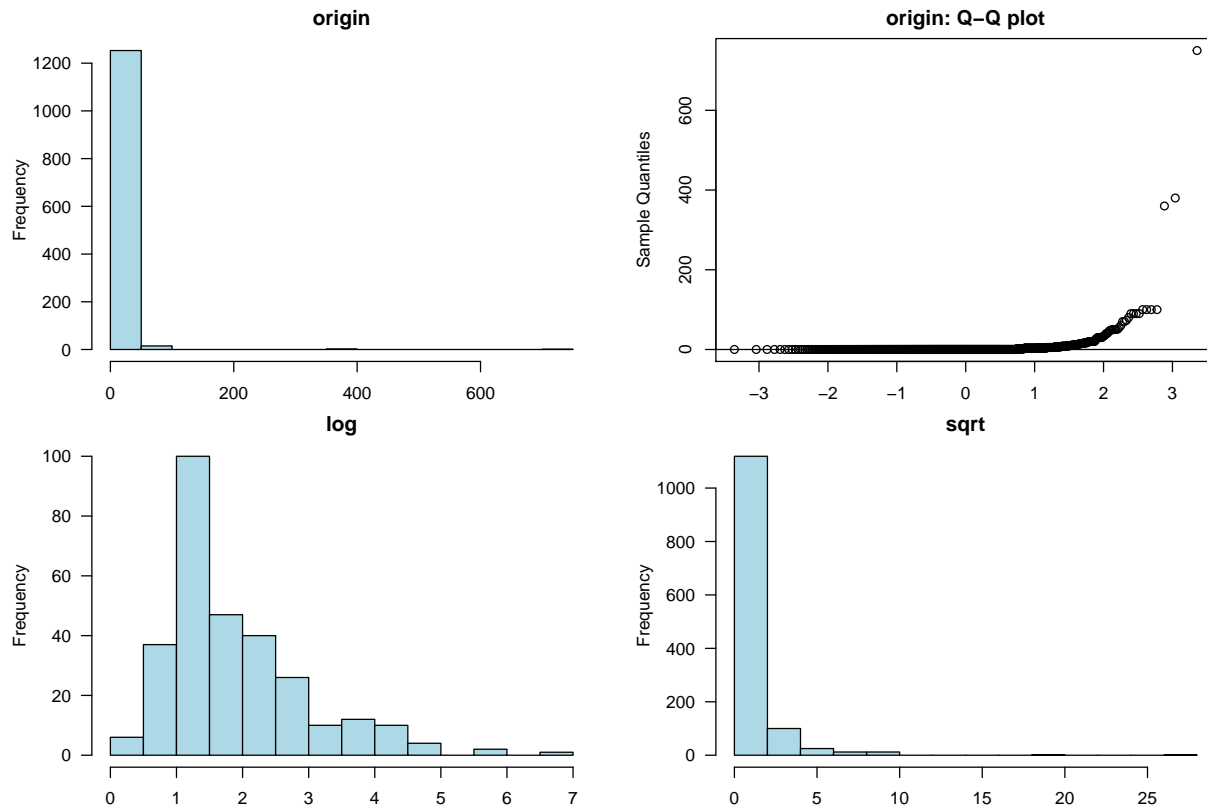
Normality Diagnosis Plot (UCB_Calculated_Invited_Attendees__c)



Normality Diagnosis Plot (UCB_Calculated_Number_of_Actual_Attendee_c)



Normality Diagnosis Plot (UCB_Actual_Attendees_c)



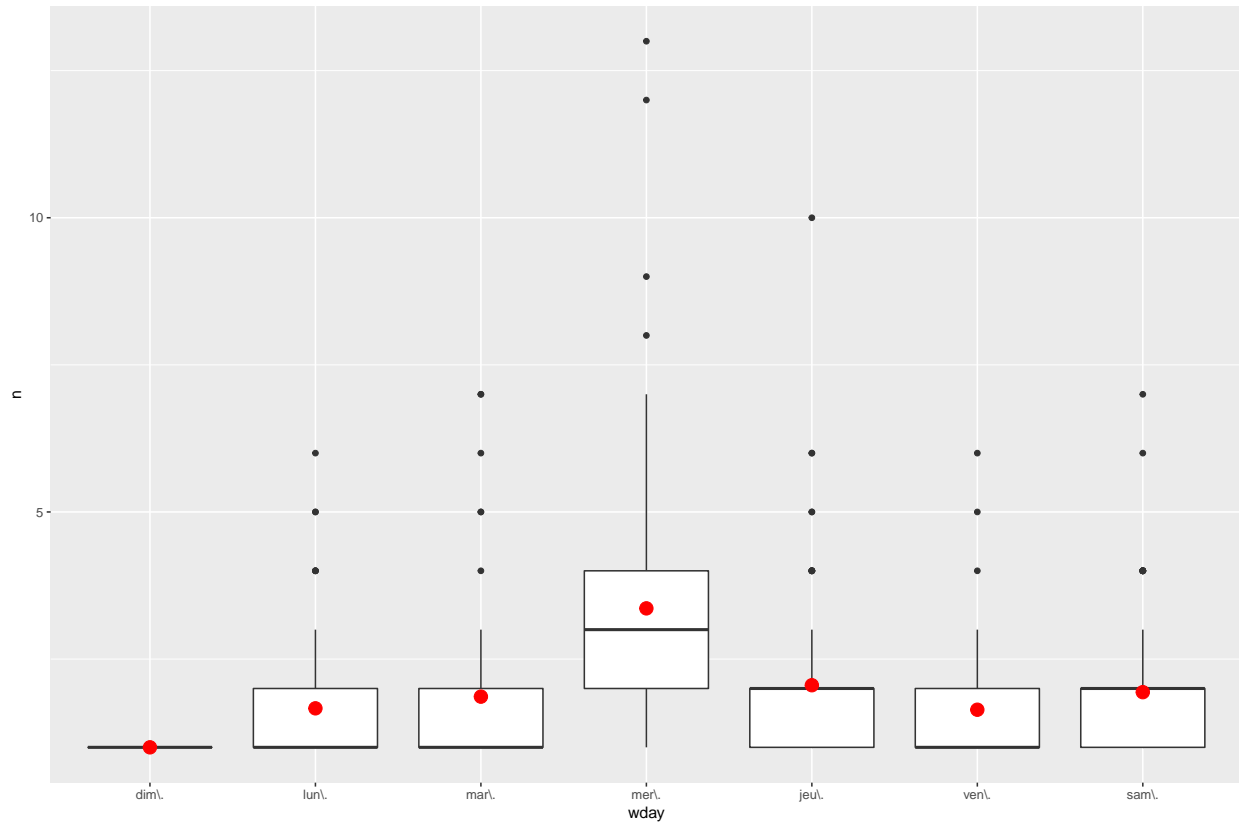
4.2 Detecting Outliers

day-of-week effect that dominates the subtler patterns we fit the model (lm) and display its predictions overlaid on the original data

Let's figure out what data are getting wrong. First, I plotted the distribution of the data, and saw that it follows a normal distribution with mean = , std = To classify a data point as an outlier, I considered any data ± 2 std above/below the mean is a significant error. Then, I created two data frames. One consisting of the "outliers", and another consisting of the "average". I plotted these on a graph (red = outliers; green = average), outlining the shape of the data. We observe something interesting.

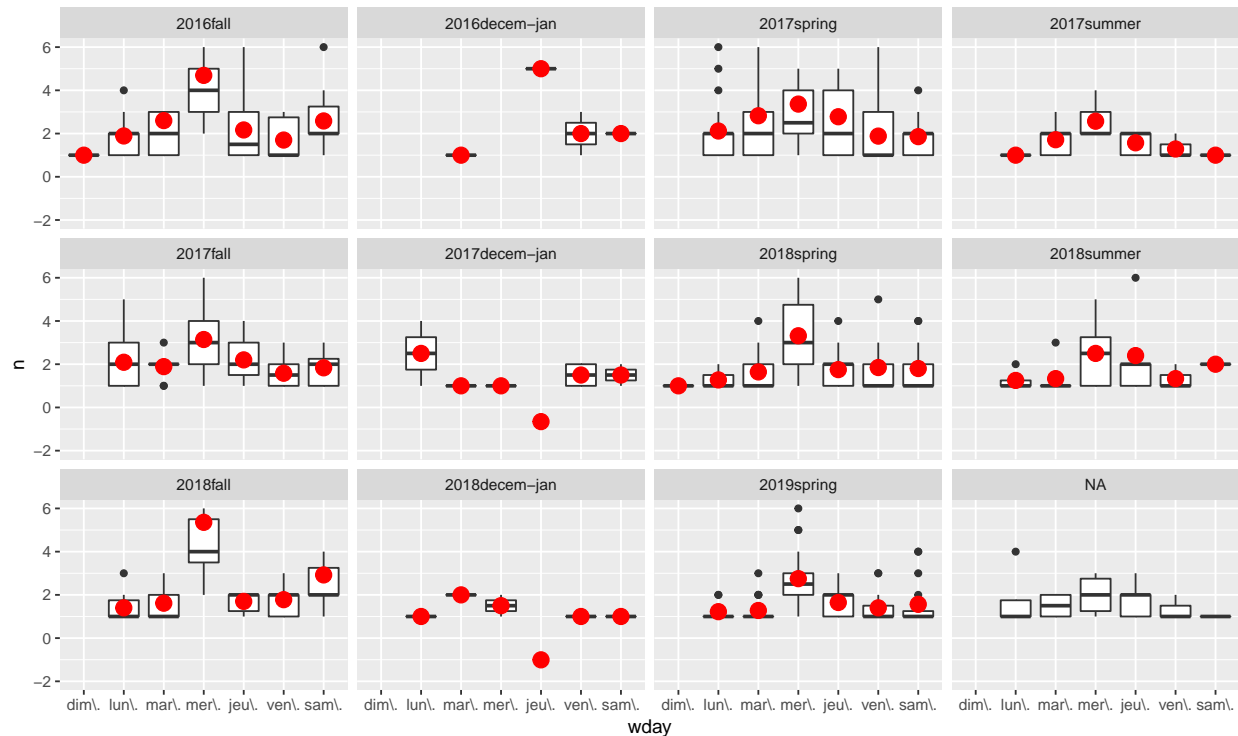
pdf

2



There is significant variation across the terms,so fitting a separate day-of-week effect for each term

pdf
2



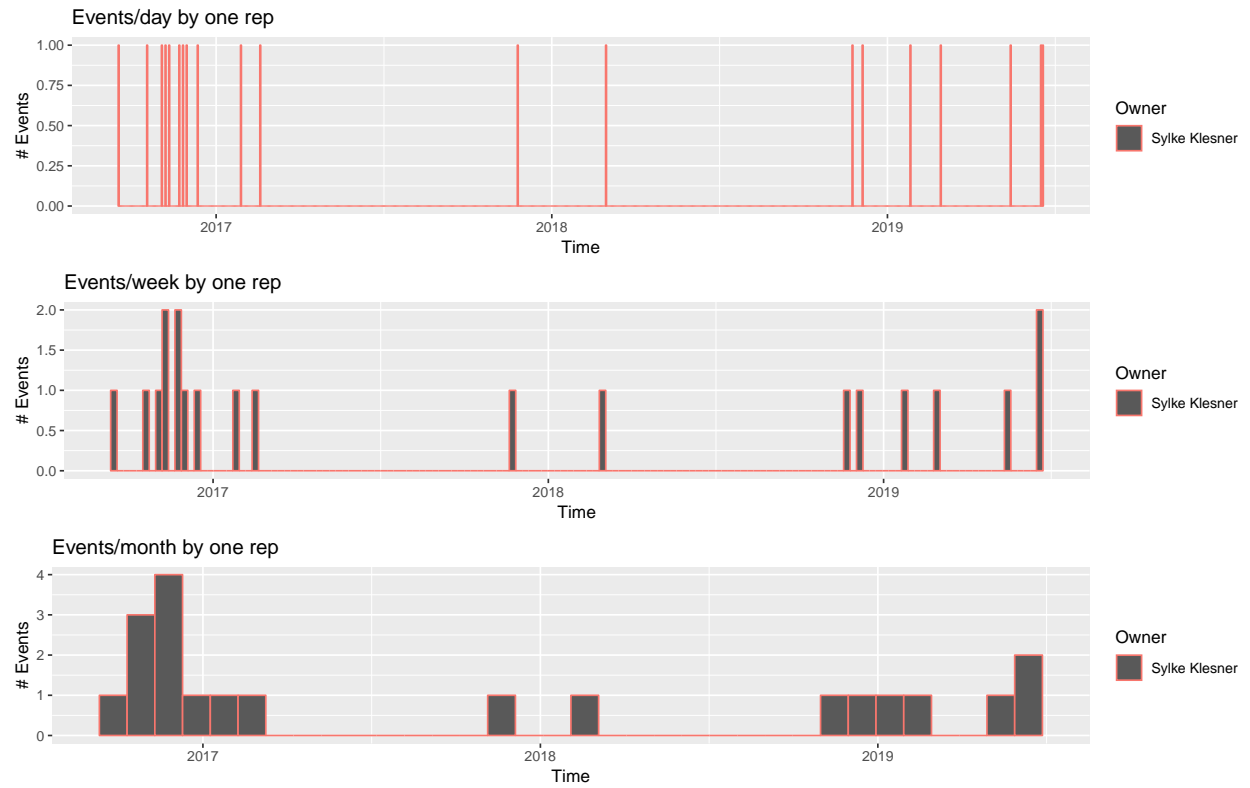
Understand how often each rep events on customers How many reps were involved in events

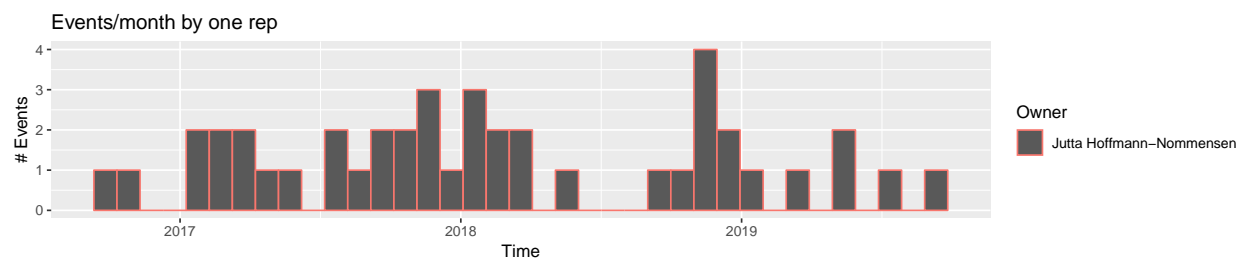
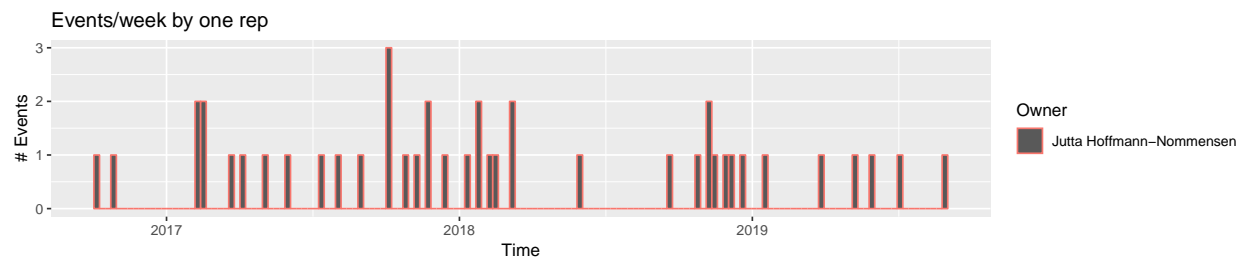
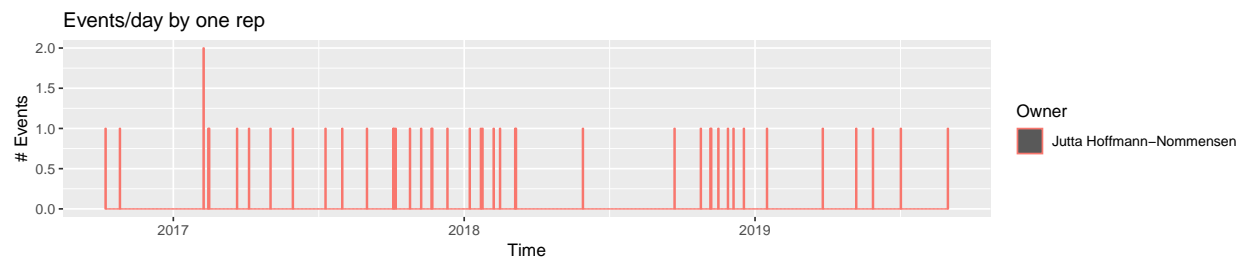
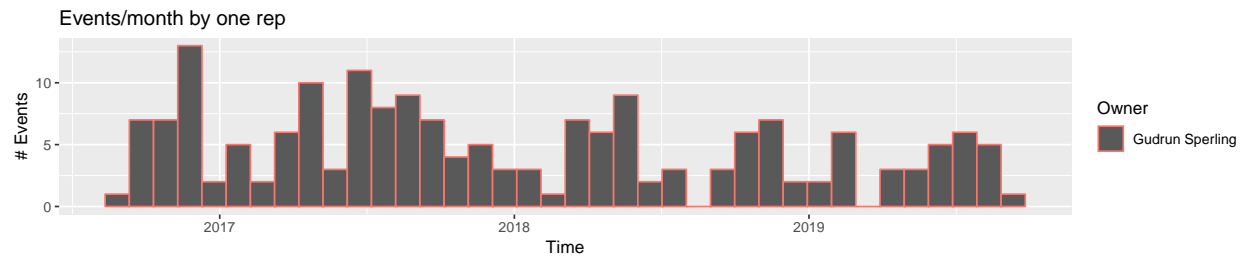
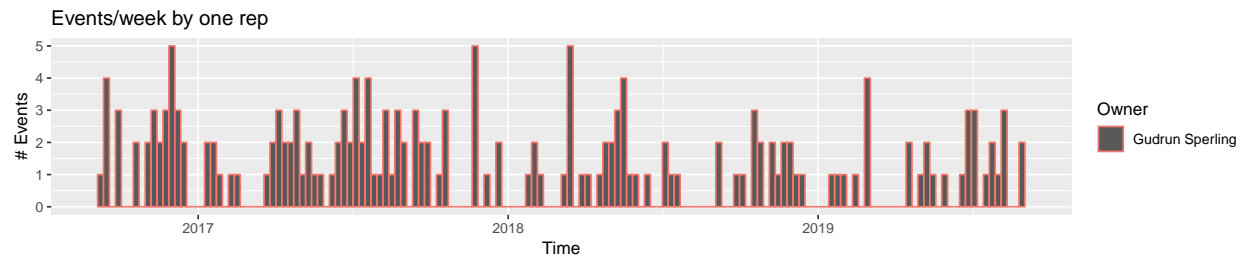
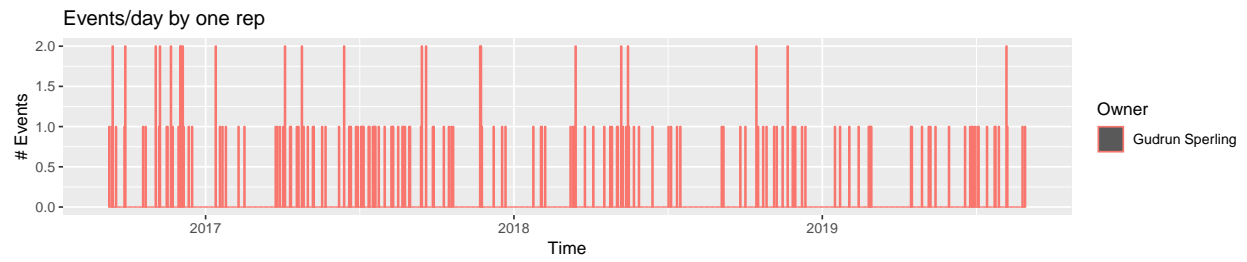
```
[1] List of reps who carried out events in 201608-201908
[1] "Sylke Klesner"           "Gudrun Sperling"
[3] "Jutta Hoffmann-Nommensen" "Astrid Hemschemeier"
[5] "Constanze Erwin"         "Susanne Meinel-Schwebke"
[7] "Nico Haag"               "Lambrini Tontsidou"
[9] "Bodo Stelle"             "Nadine von Ryssel"
[11] "Friederike Hilger"       "Petra Opolka"
[13] "Angelika Gundlach"       "Sarah Sahl"
[15] "Udo Lendl"               "Stephanie Mueller-Varain"
[17] "Eike-Hans Zimmermann"    "Katrin Riedel"
[19] "Klaus Kalcum"            "Rainer Krüger"
[21] "Balint Szilagyi"         "Petra Rösler"
[23] "Angela Hertel"           "Lubica Vetrikova"
[25] "Ulrike Hermsdörfer-Vonalt" "Alya Kokot"
[27] "Renate Bohnsack"         "Oliver Feil"
[29] "Jan Peetz"               "Silvia Slazenger"
[31] "Caroline-Mascha Reifferscheid" "Sandra Mogk"
[33] "Jan Bretz"               "Caroline Lehnart-Betz"
[35] "Heike Holländer"         "Sandra Schink"
[37] "Daniel Bura"             "Heiko Sic"
[39] "Grit Brußig"             "Michael Kemper"
[41] "Bianca Morgenstern"      "Viola Wernik"
[43] "Jan Helfrich"            "Saskia Winter"
[45] "Fabian Paul"             "Annina Getz"
[47] "Manh Dan Nguyen"         "Steffen Hiller"
[49] "Rüdiger Thees"           "Florian Hinze"
[51] "Adrian Zander"           "Tabea Thomas"
[53] "Wojciech Dombrowsky"     "Missanga van de Sand"
```

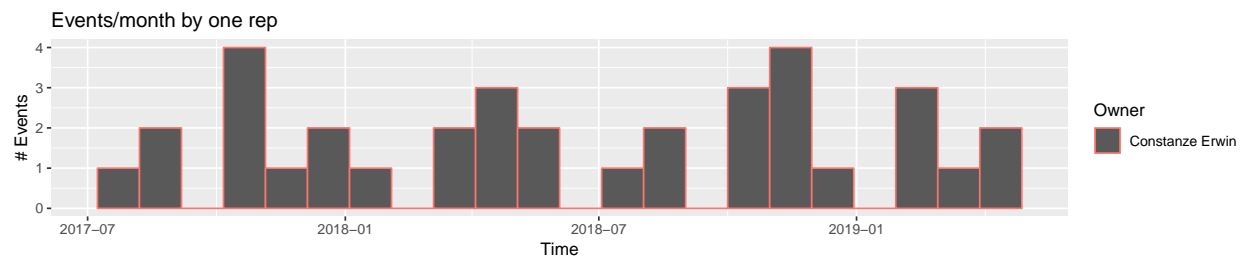
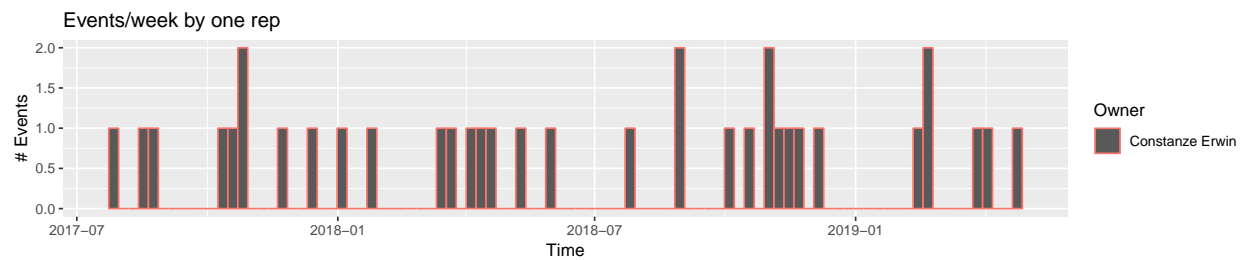
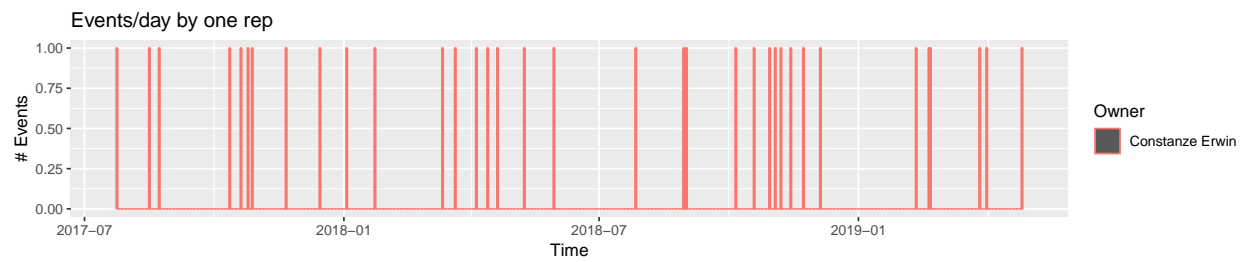
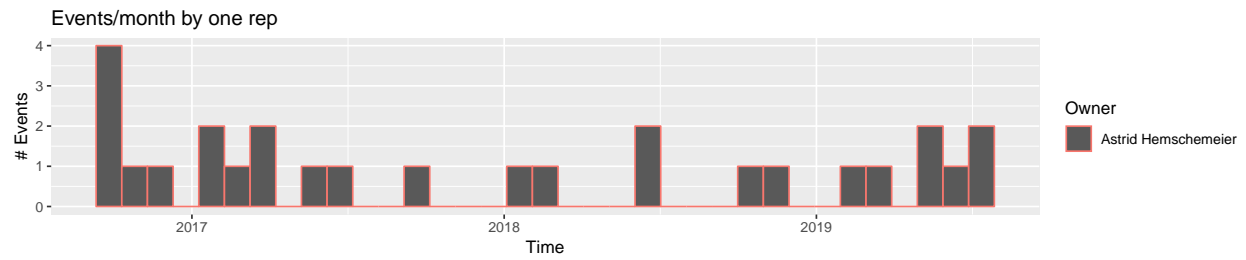
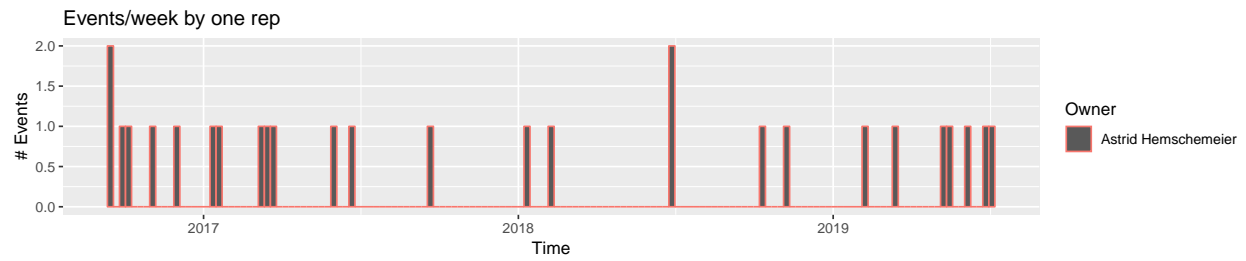
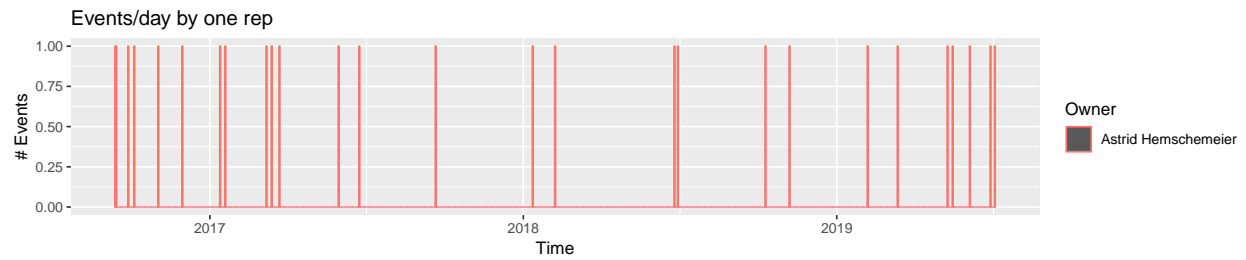
[55] "Joachim Hipp"

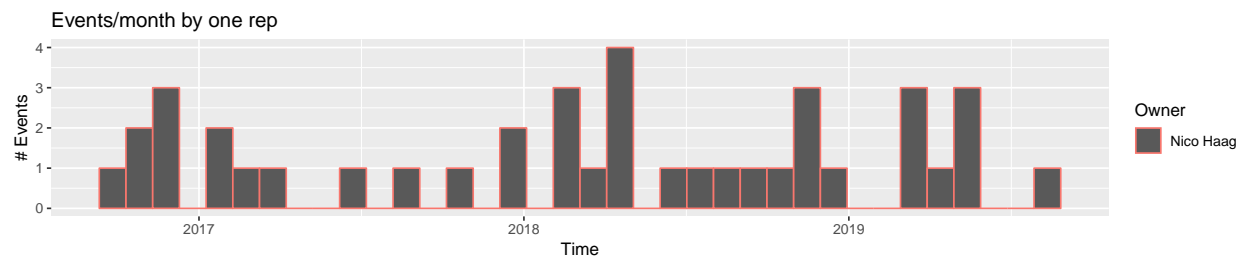
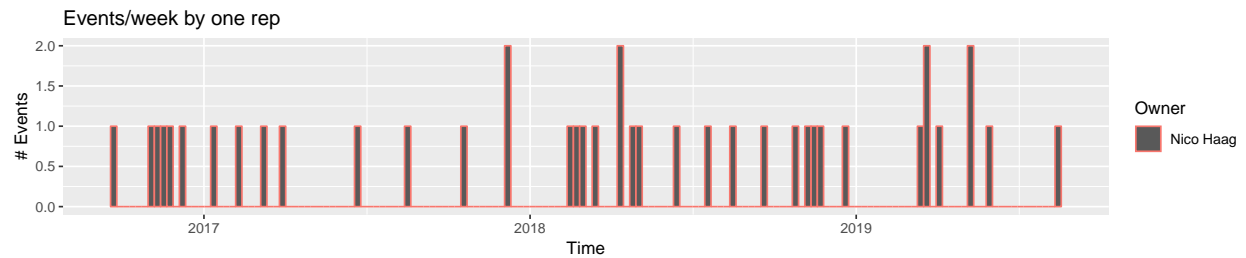
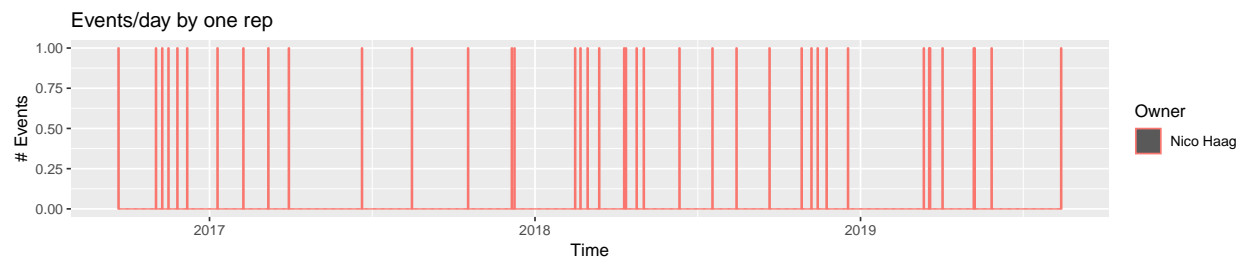
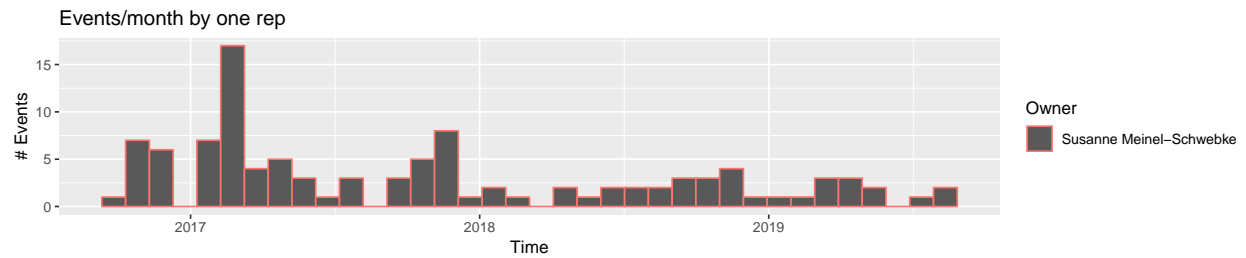
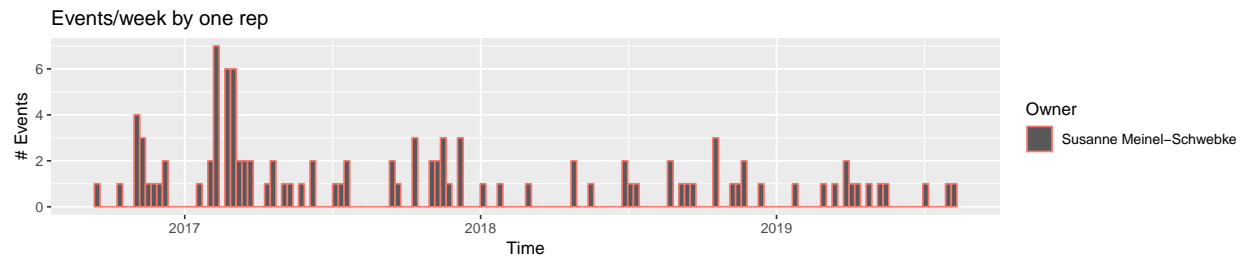
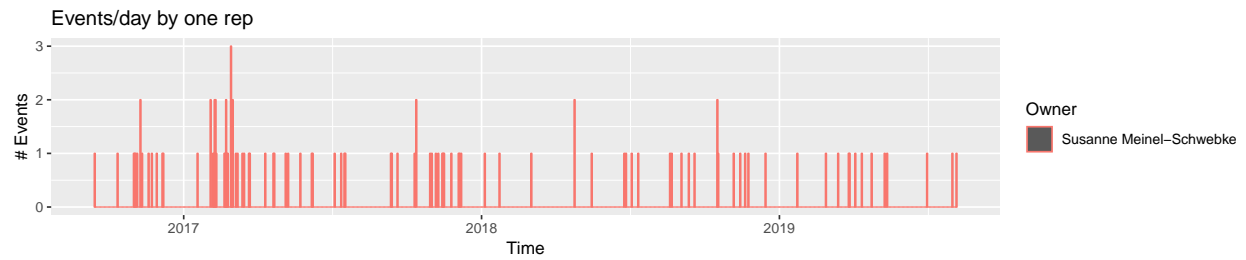
pdf

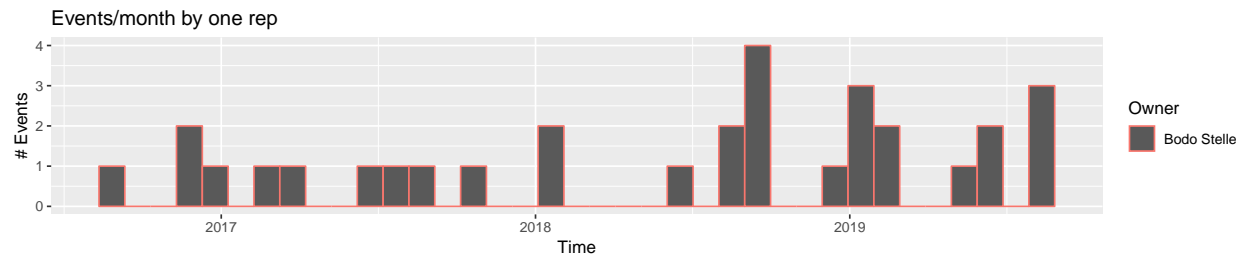
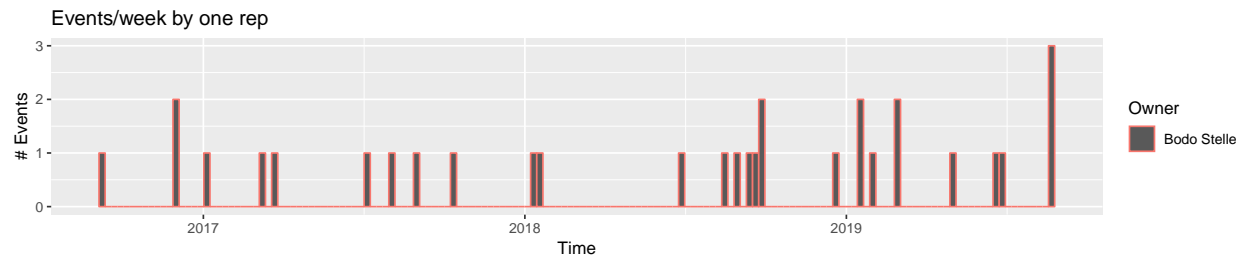
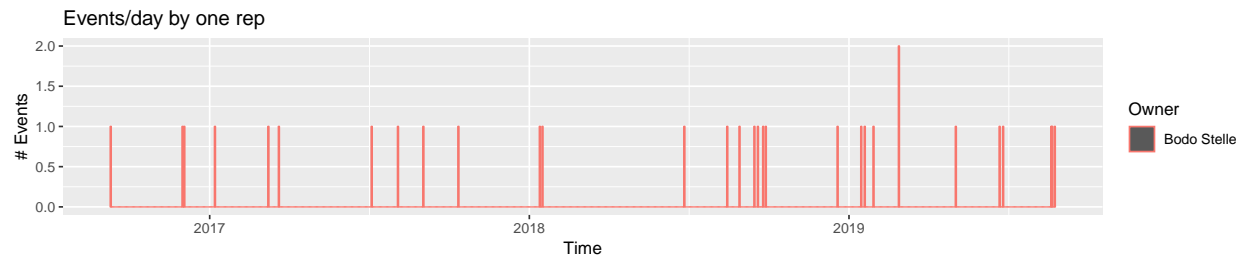
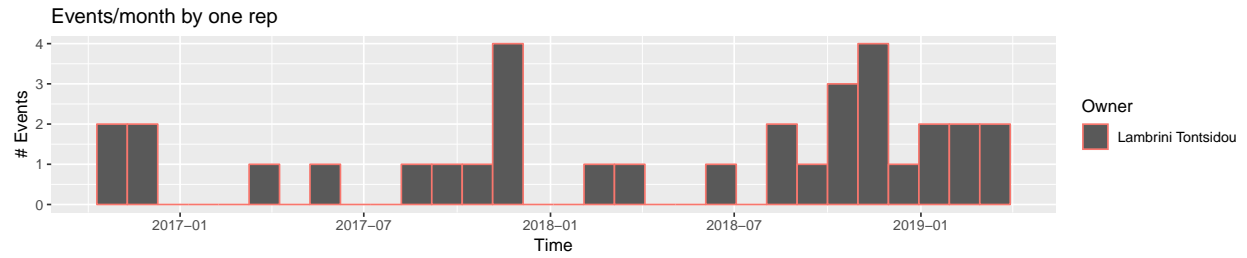
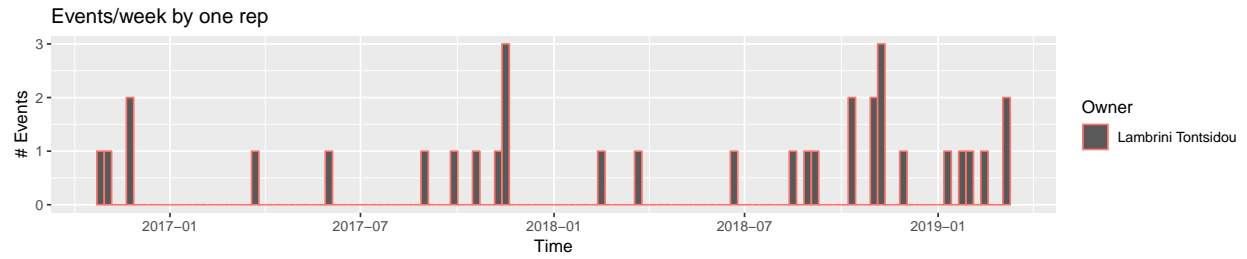
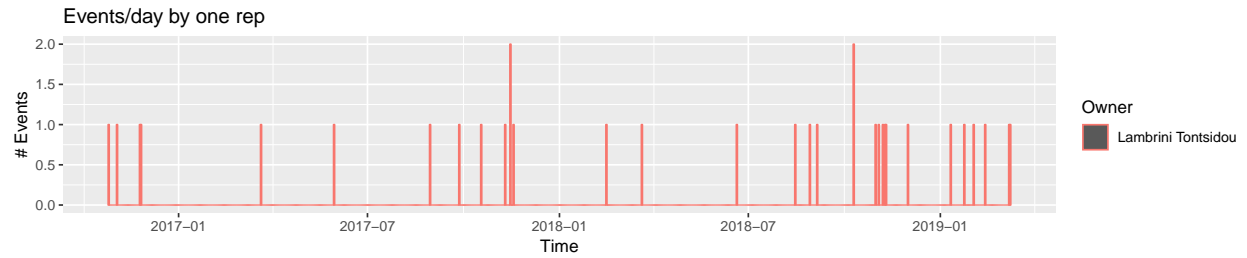
2

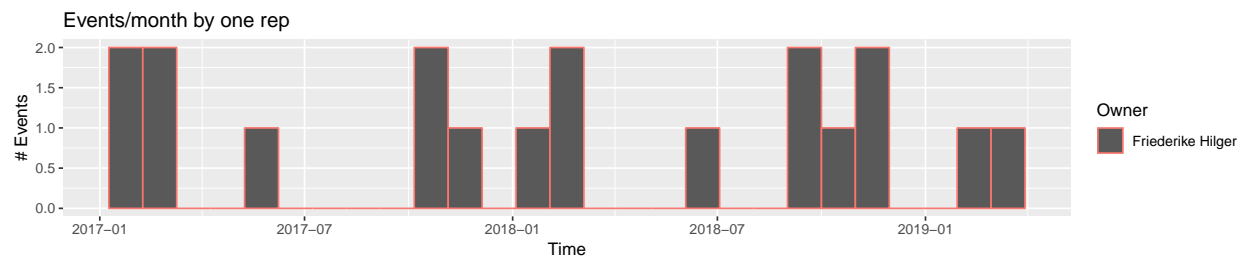
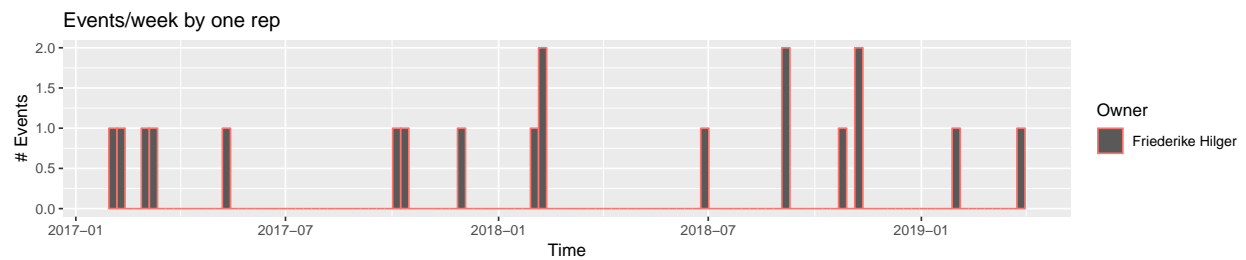
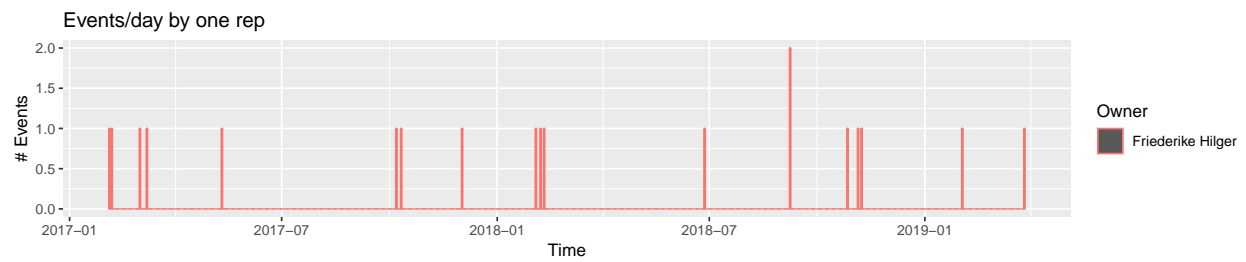
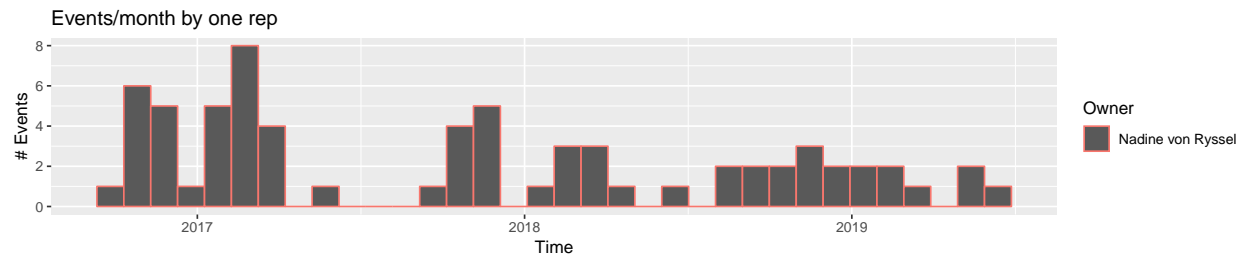
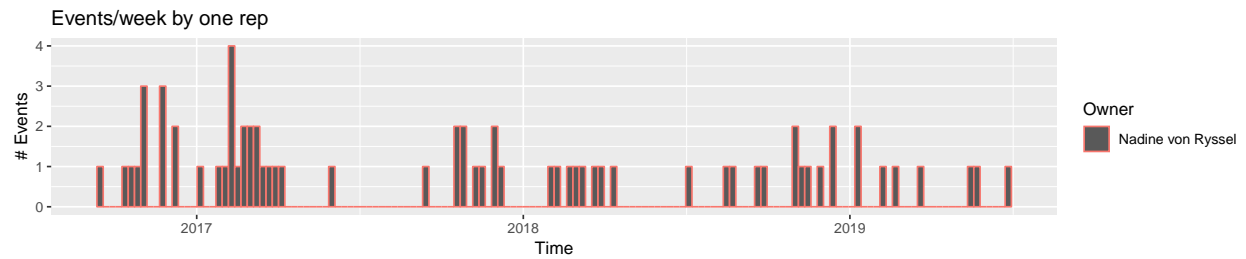
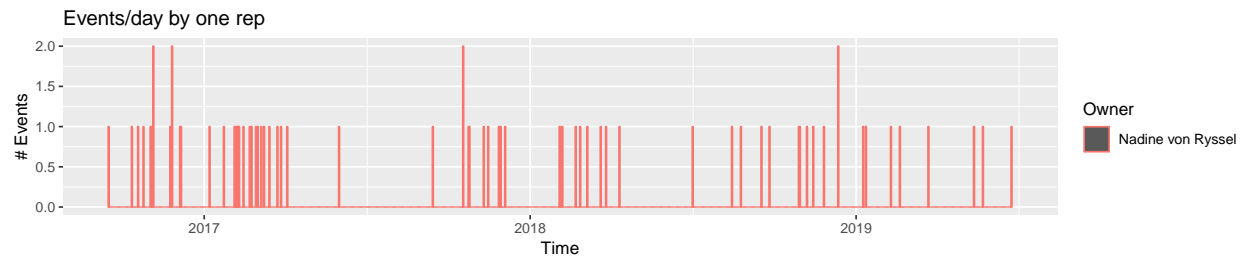


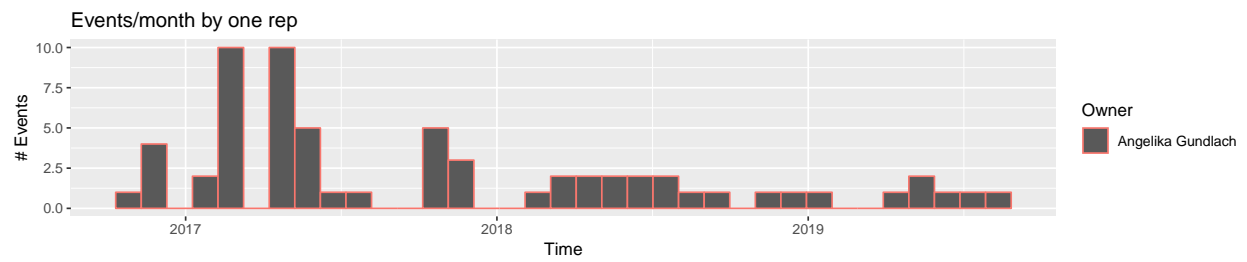
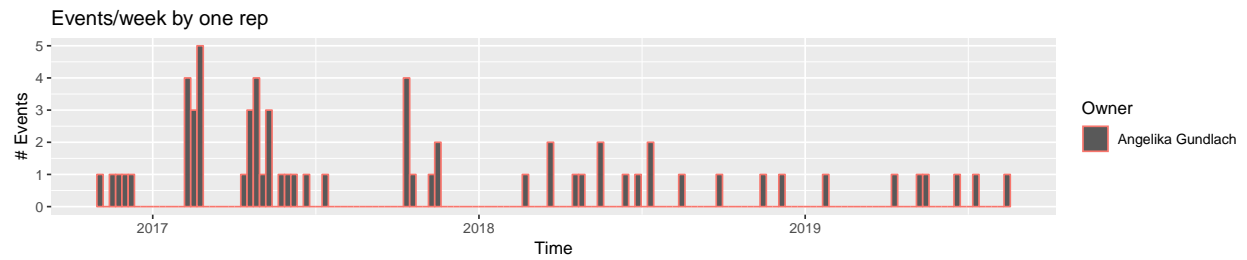
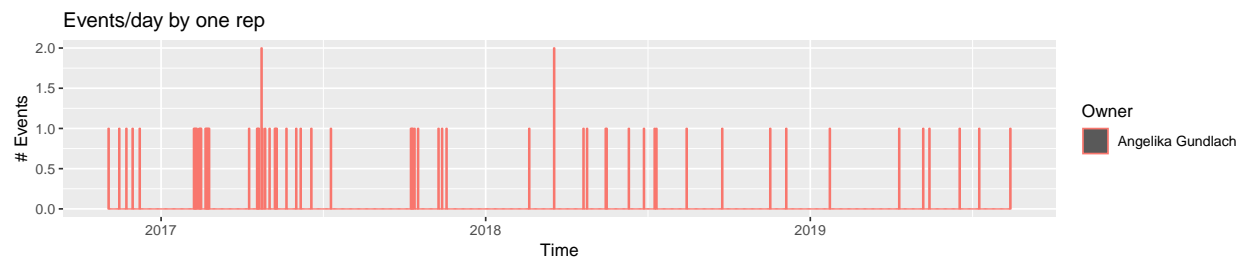
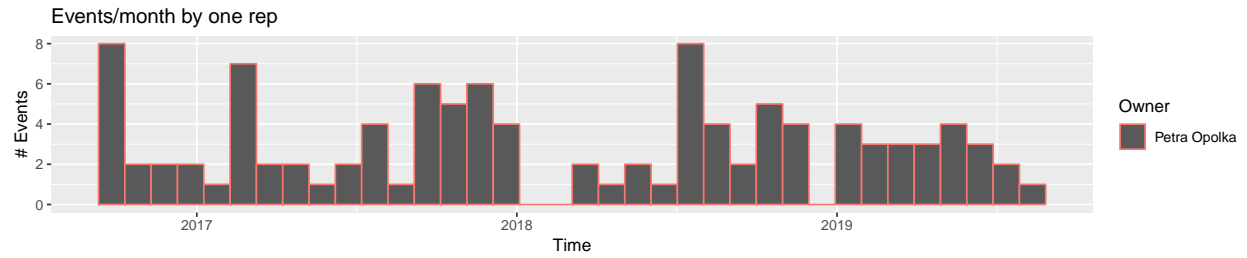
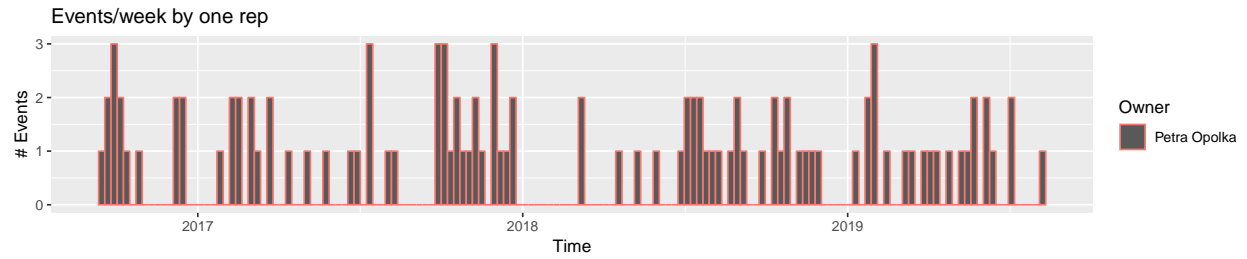
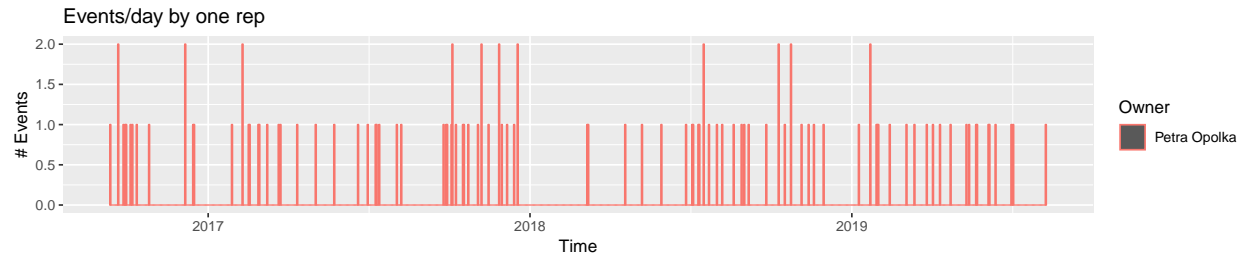


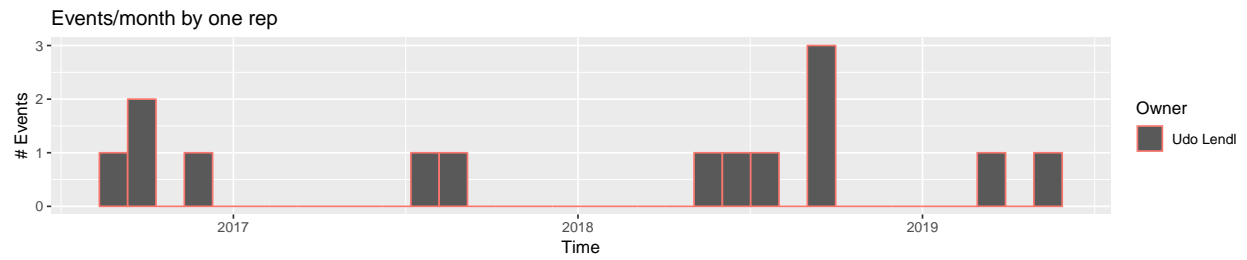
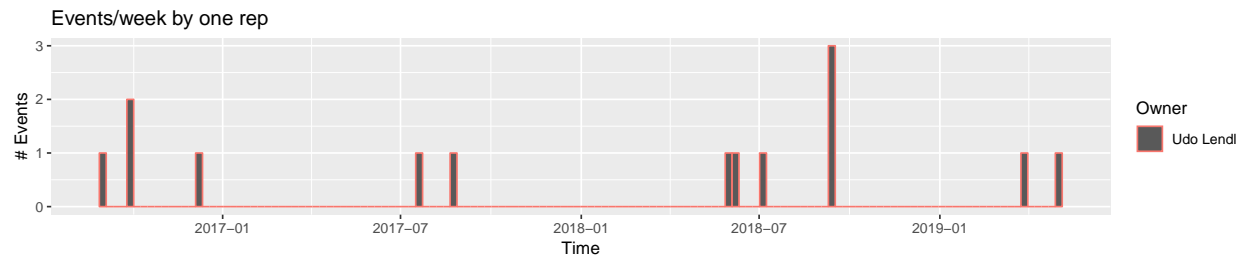
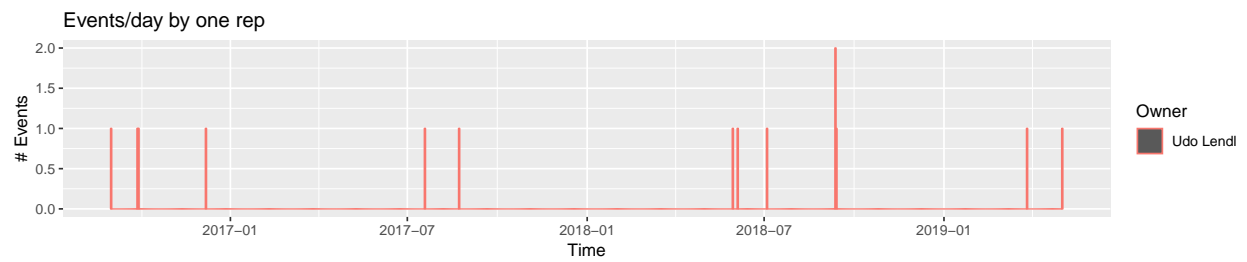
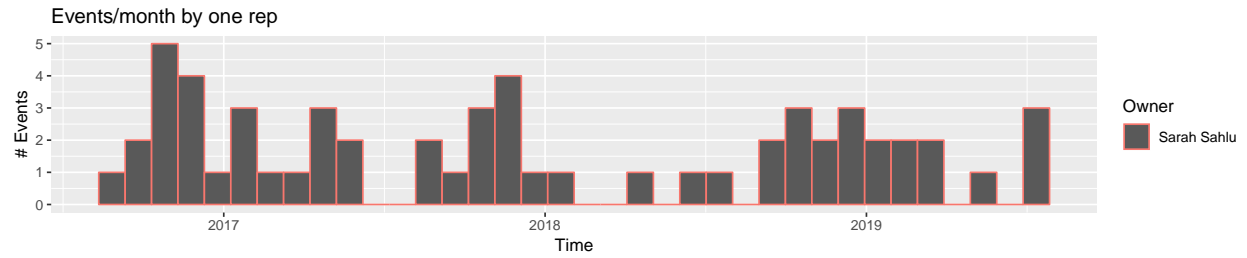
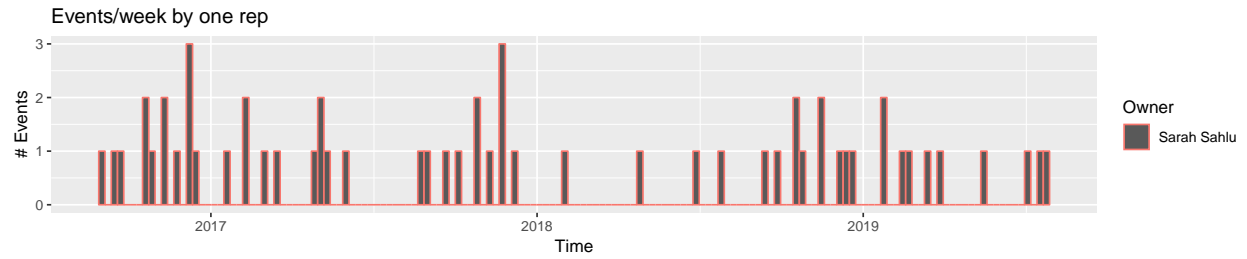
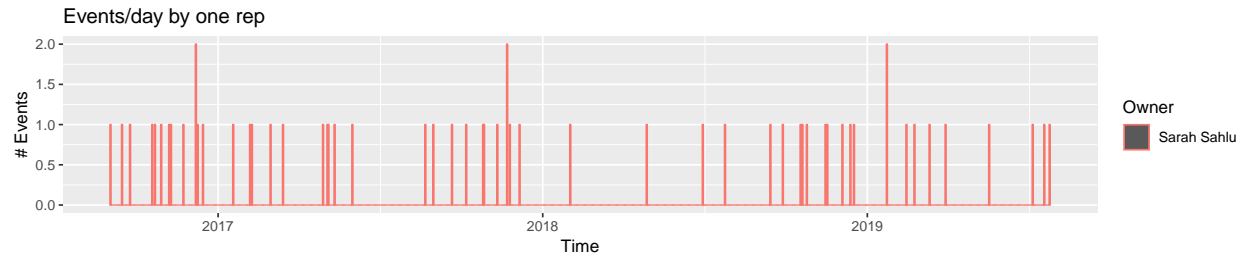


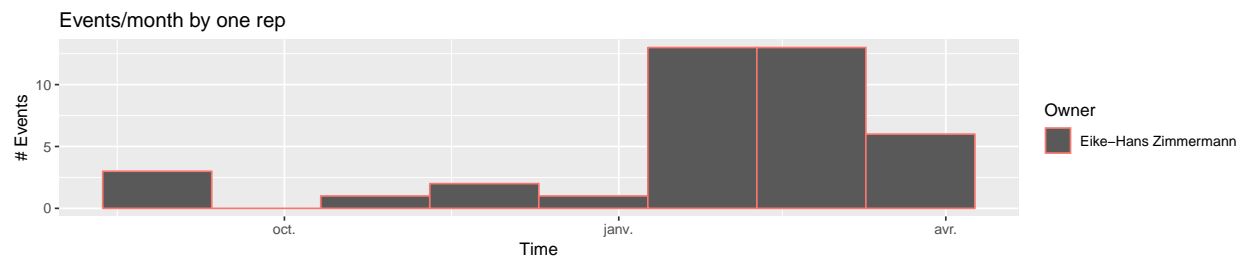
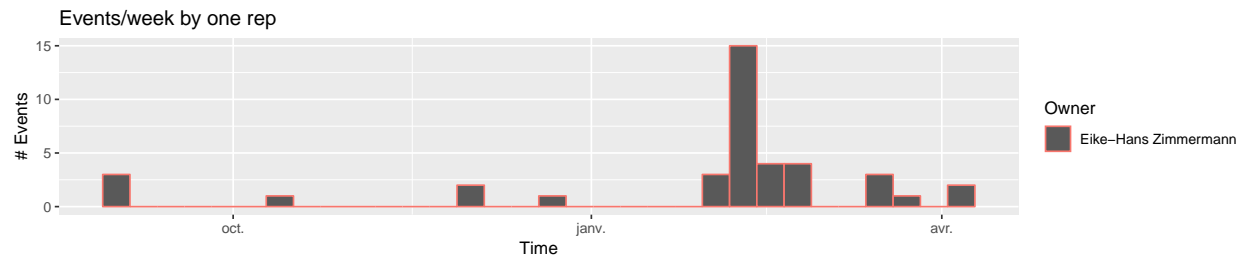
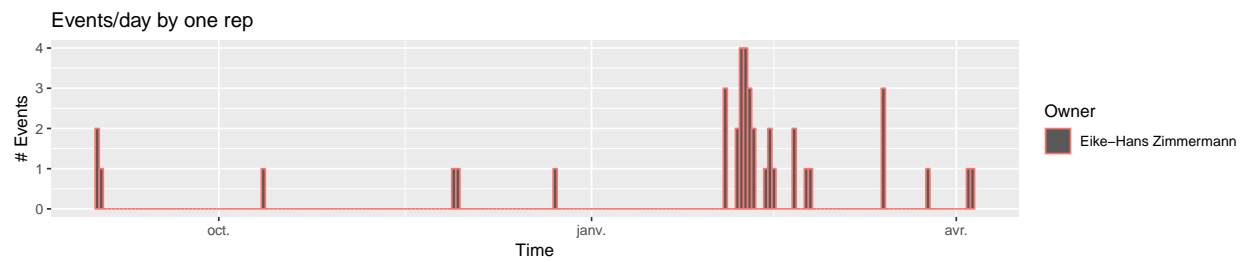
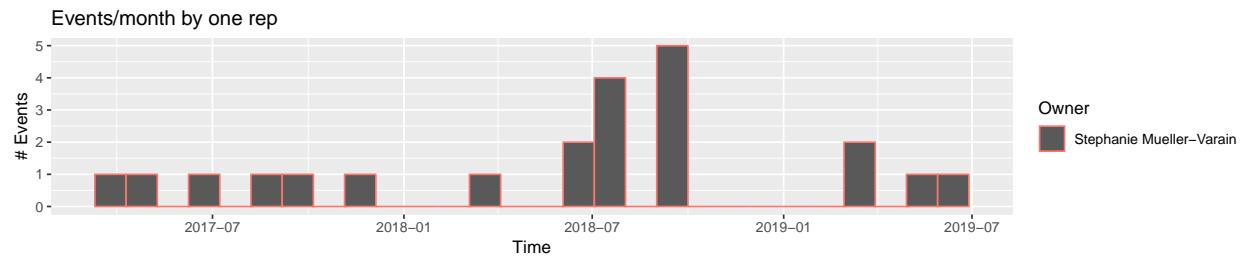
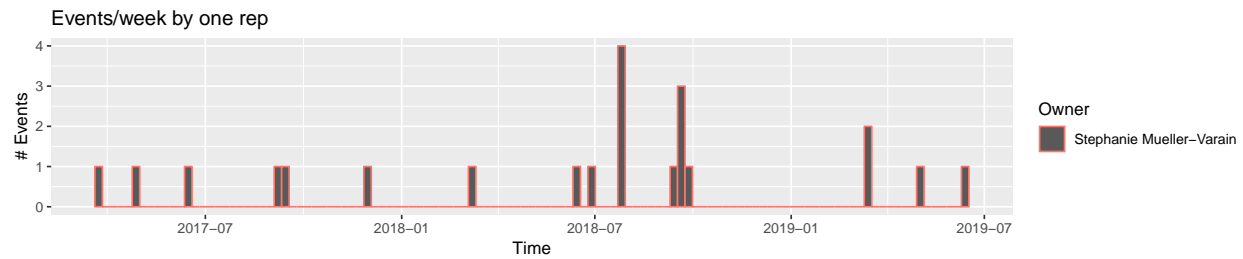
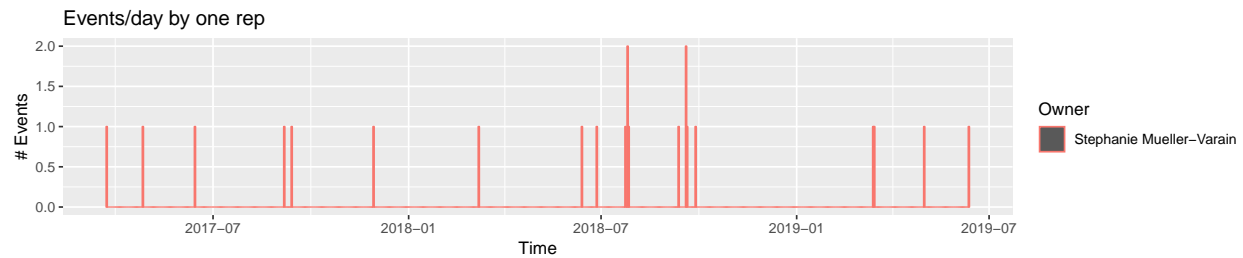


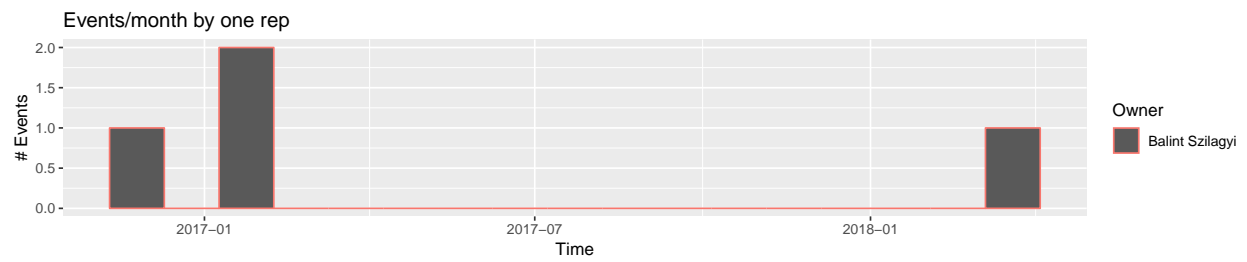
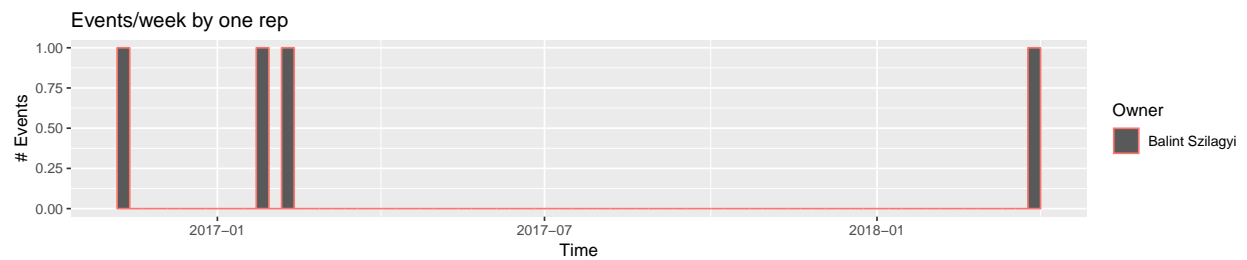
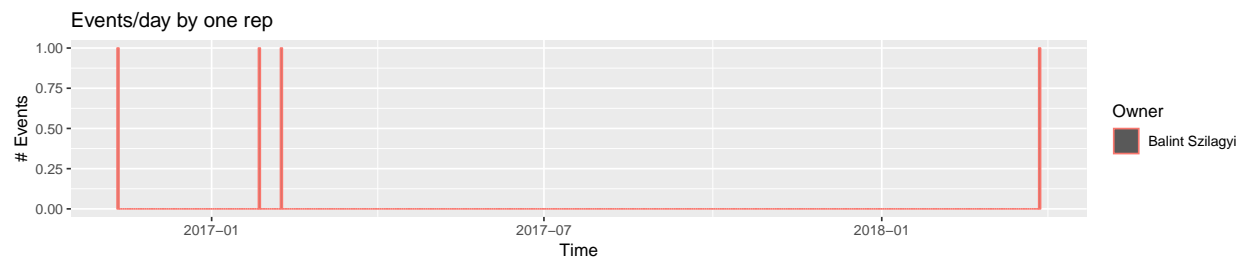
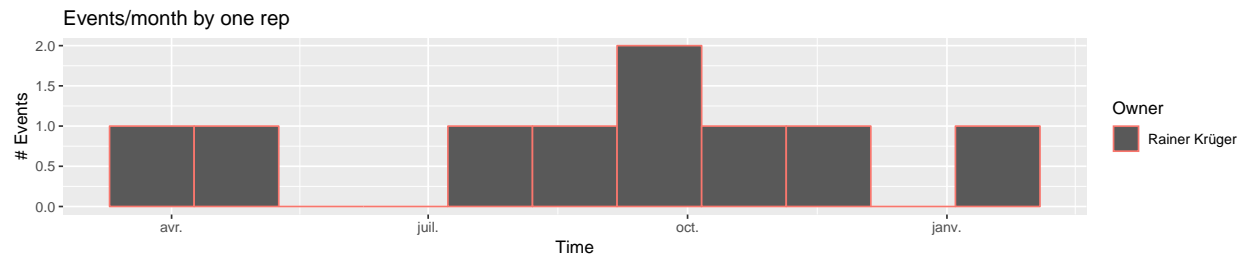
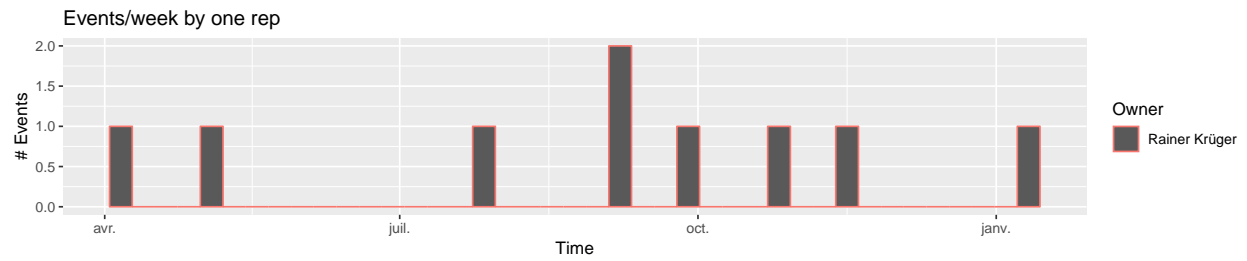
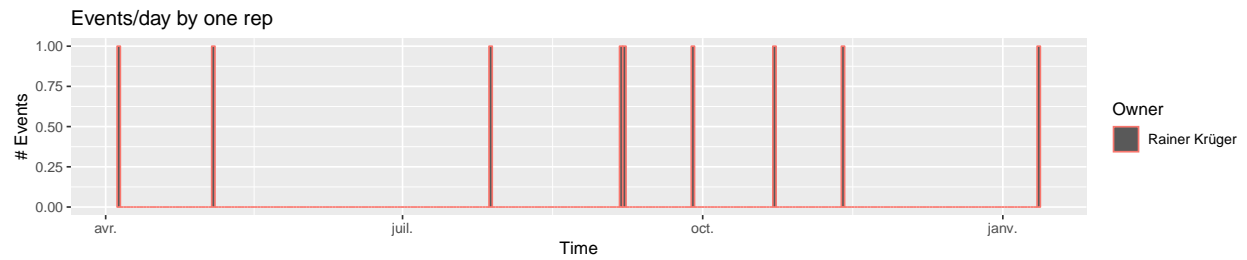


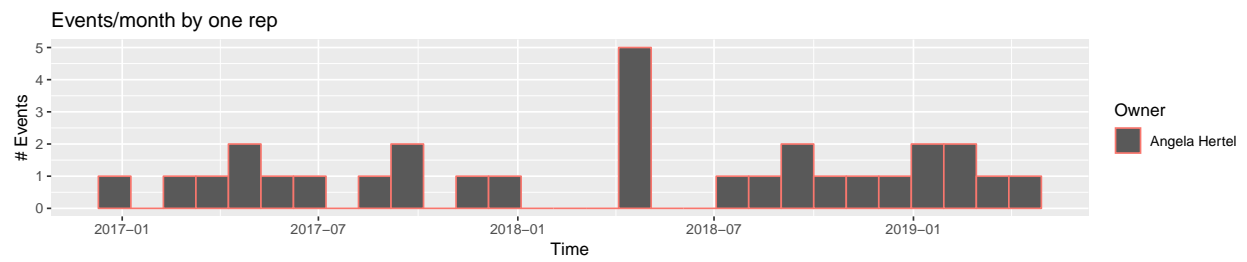
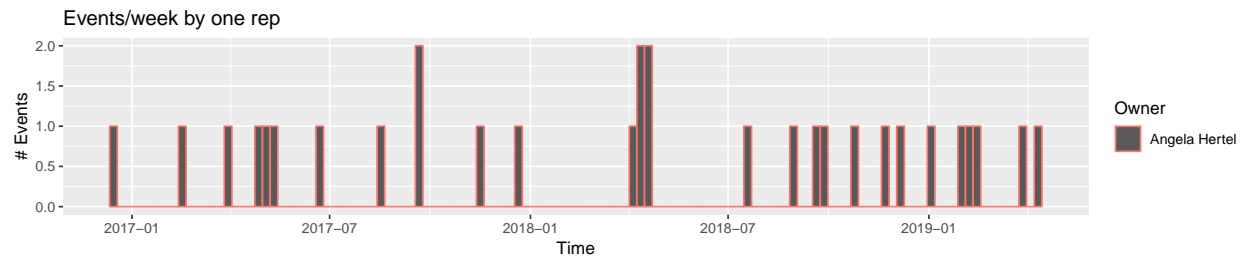
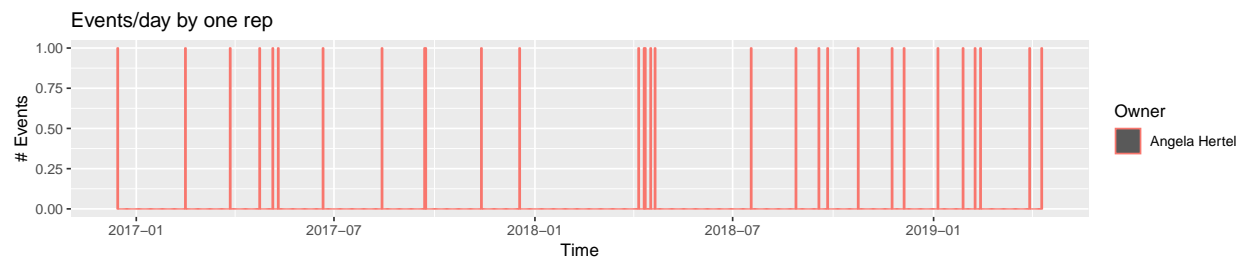
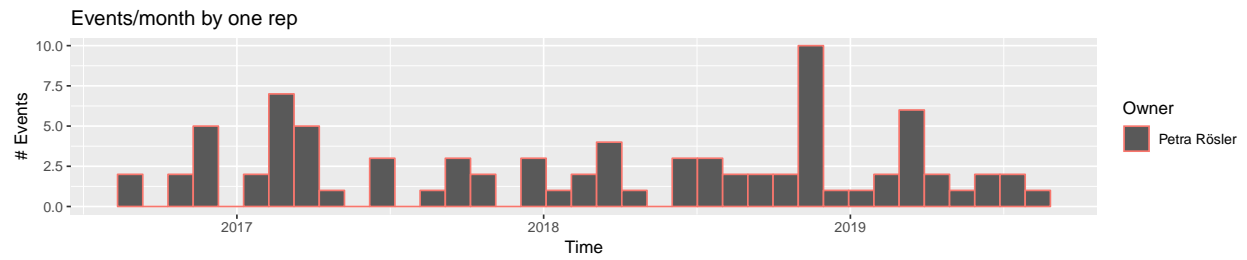
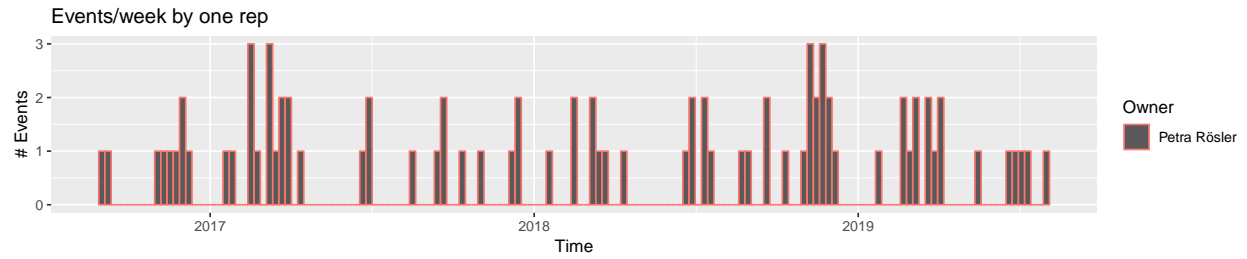
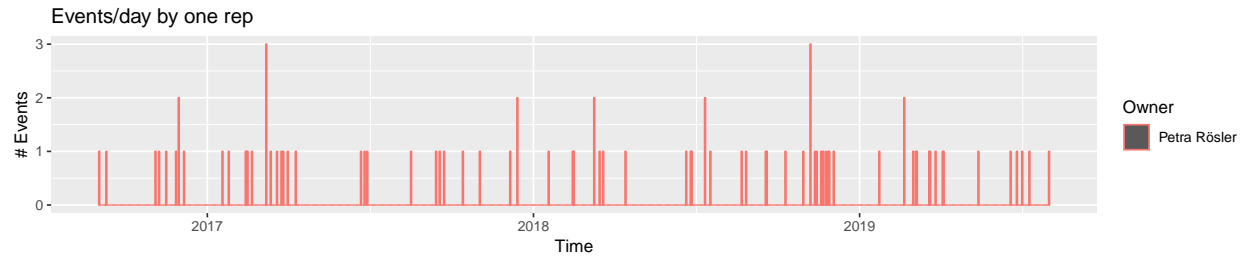


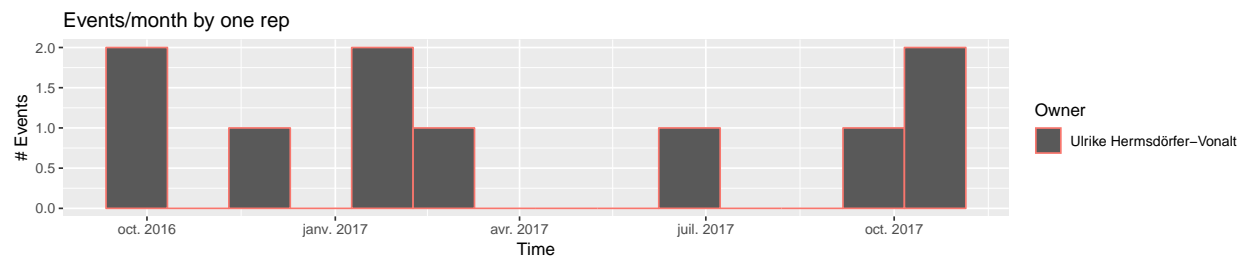
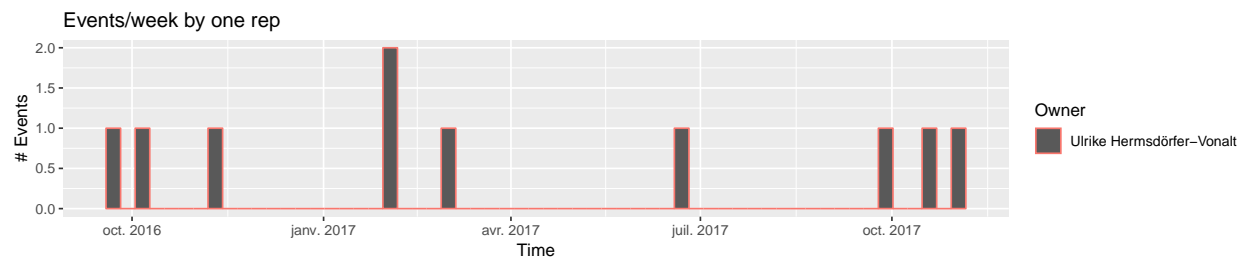
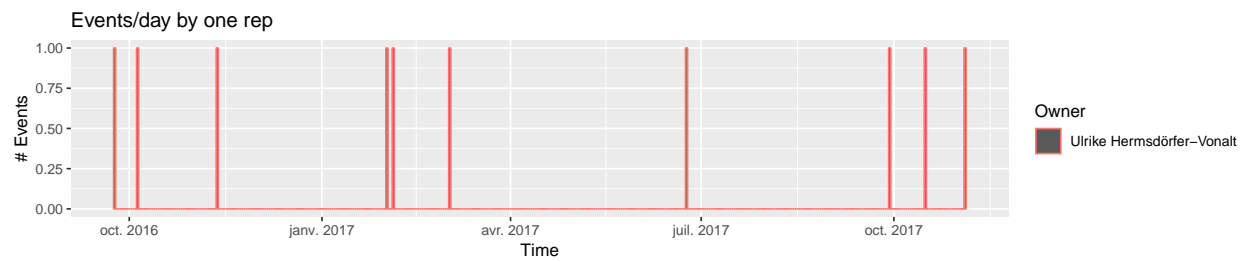
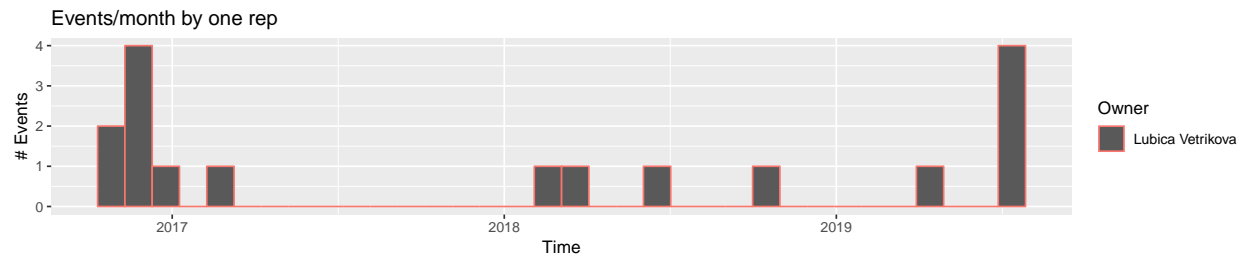
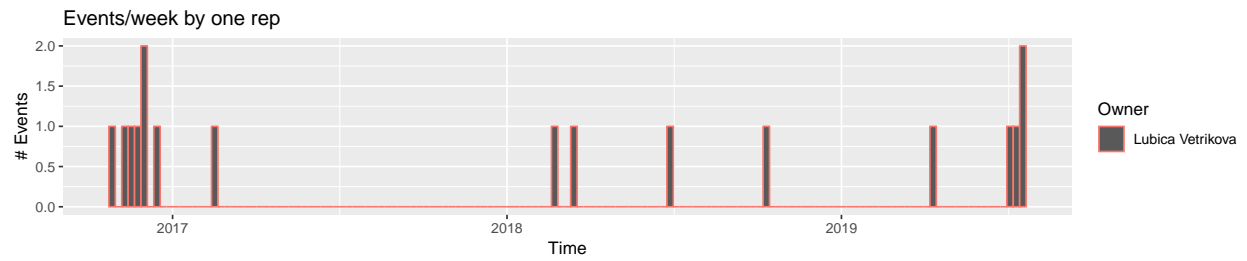
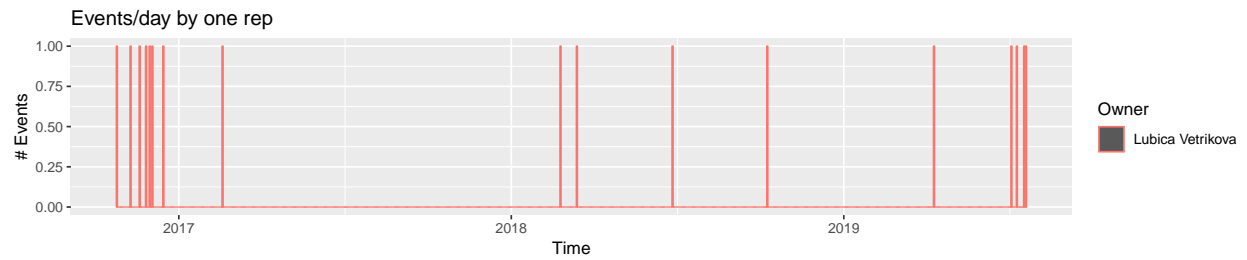


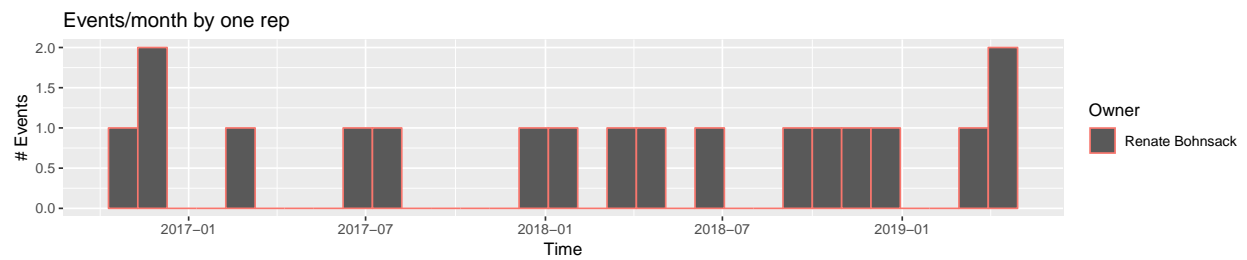
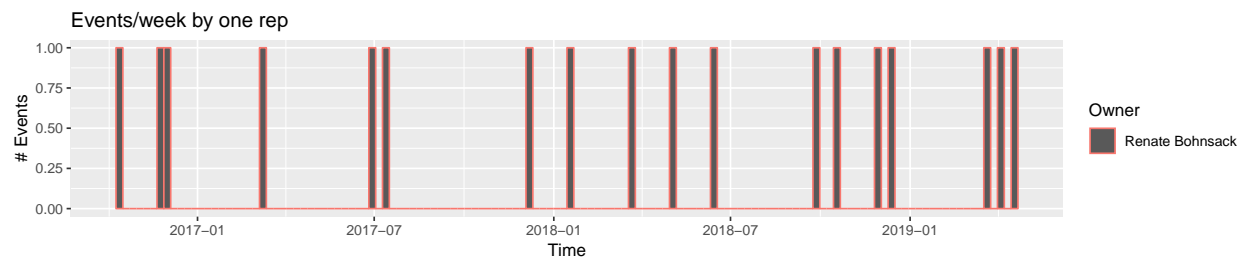
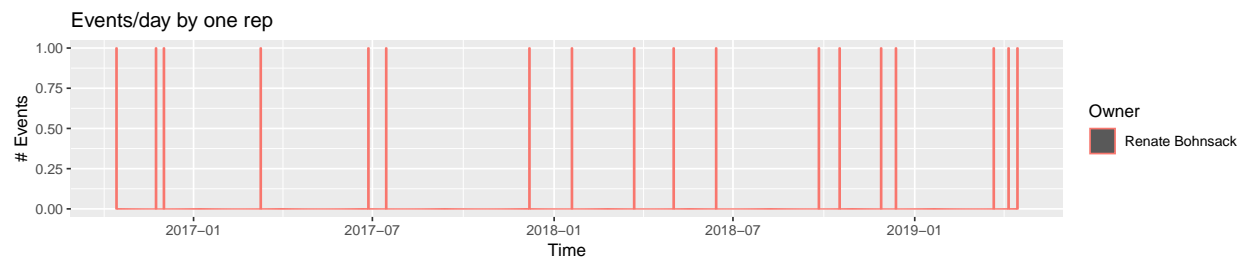
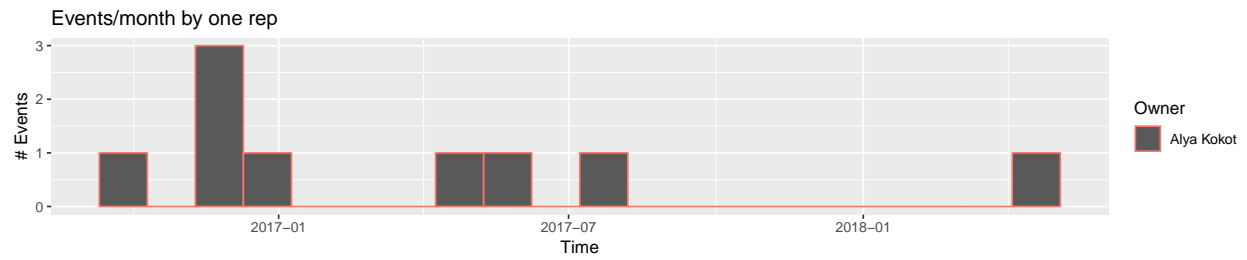
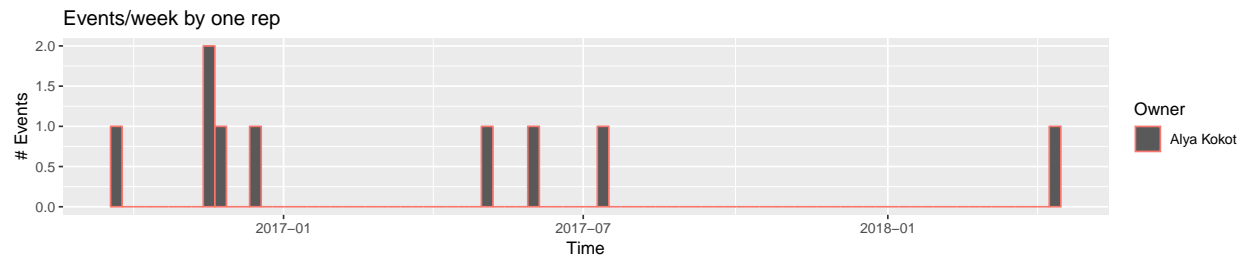
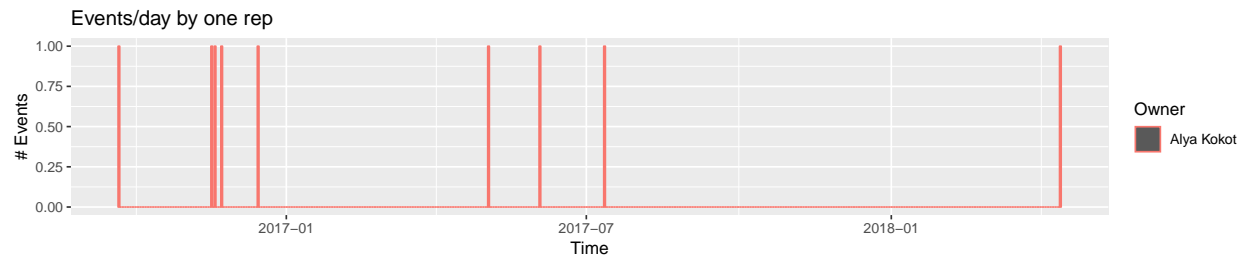


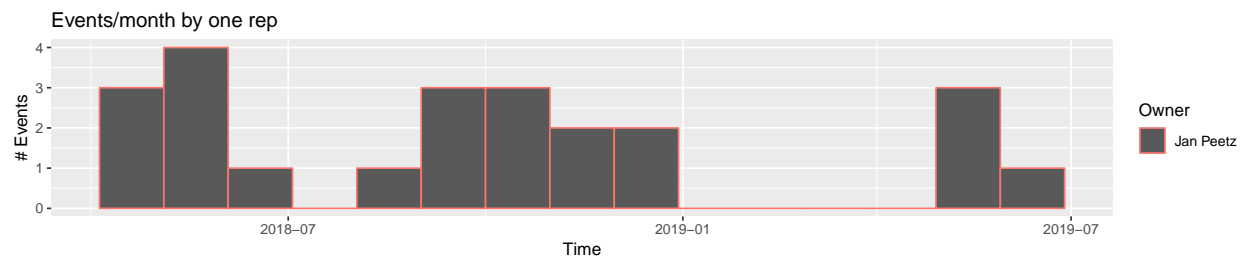
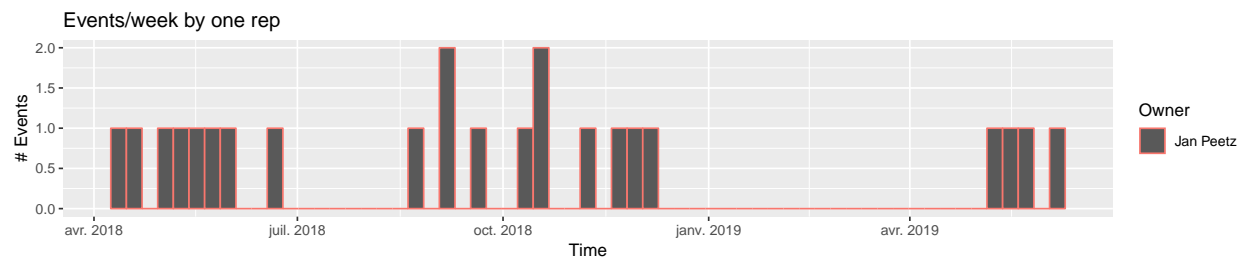
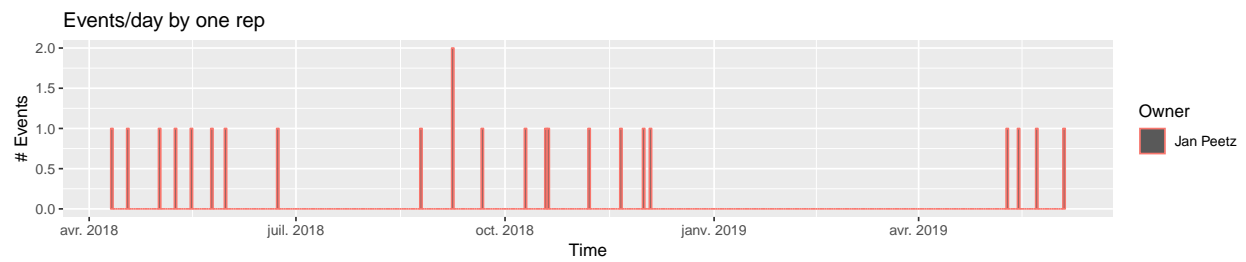
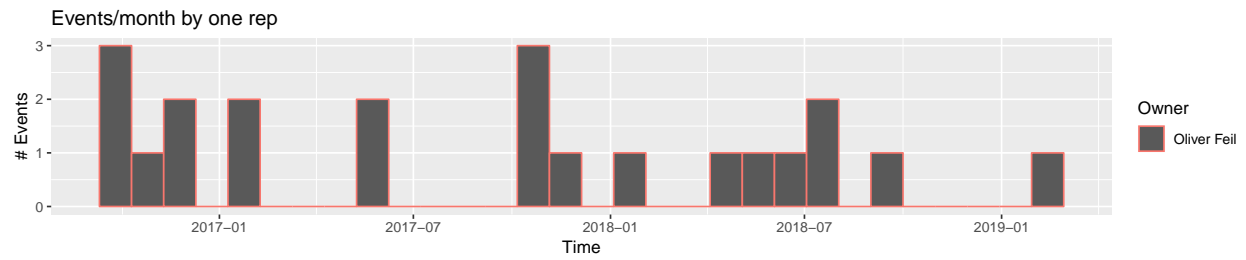
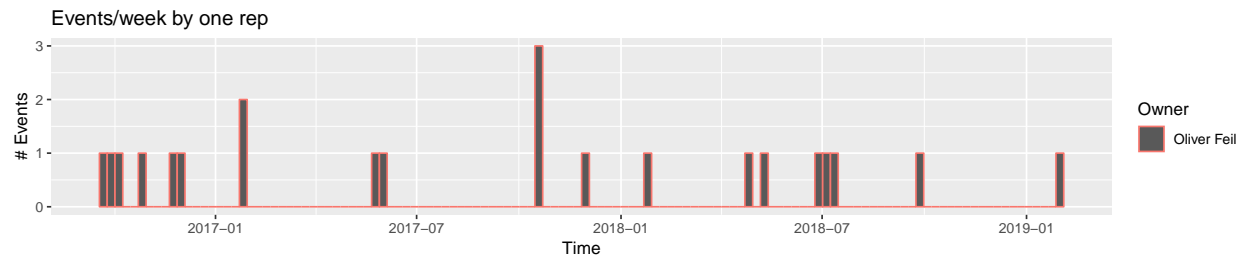
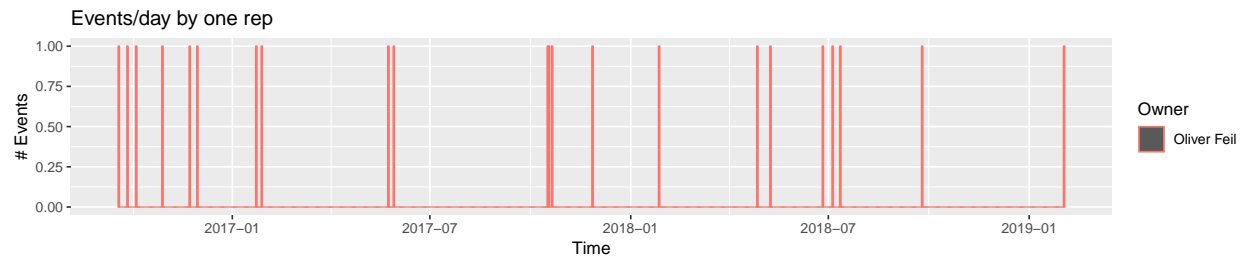


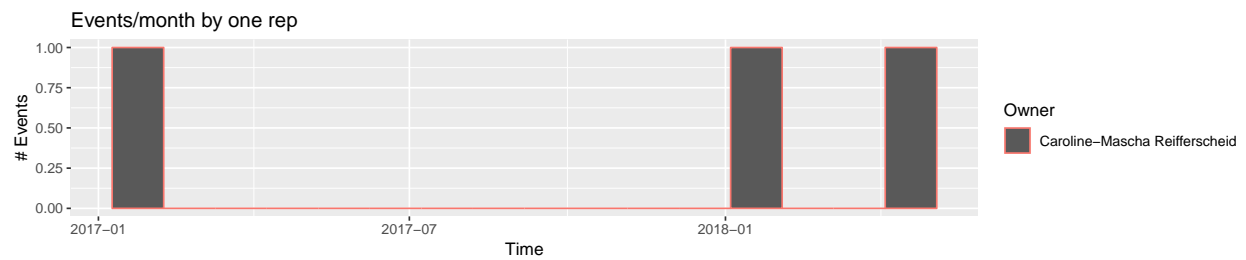
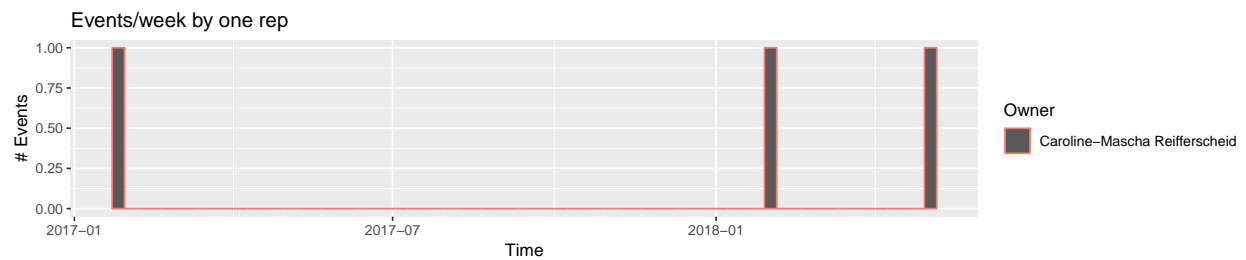
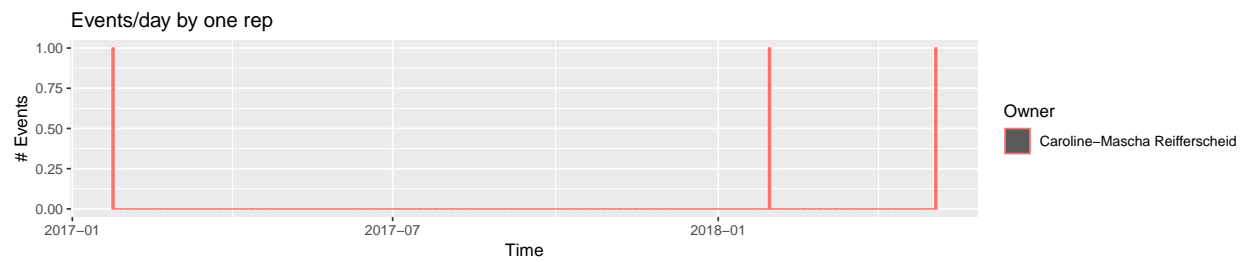
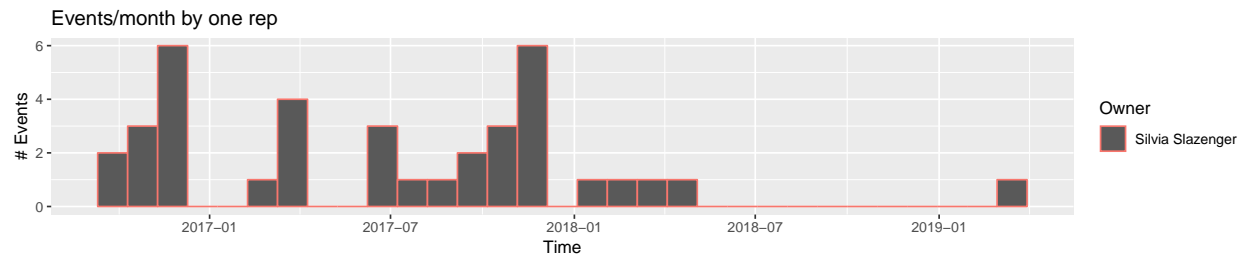
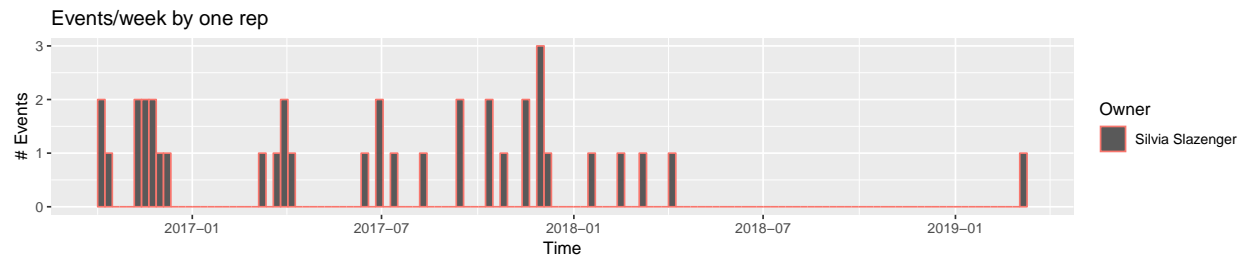
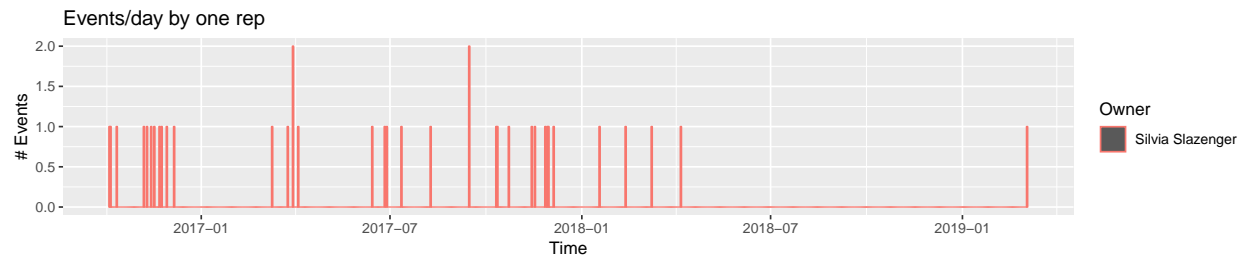


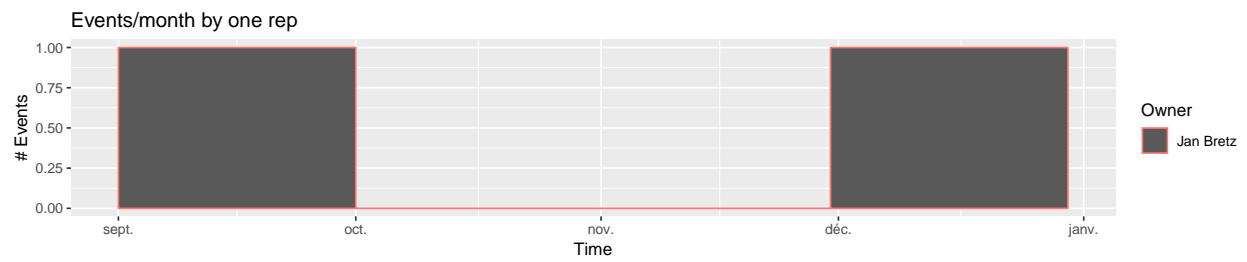
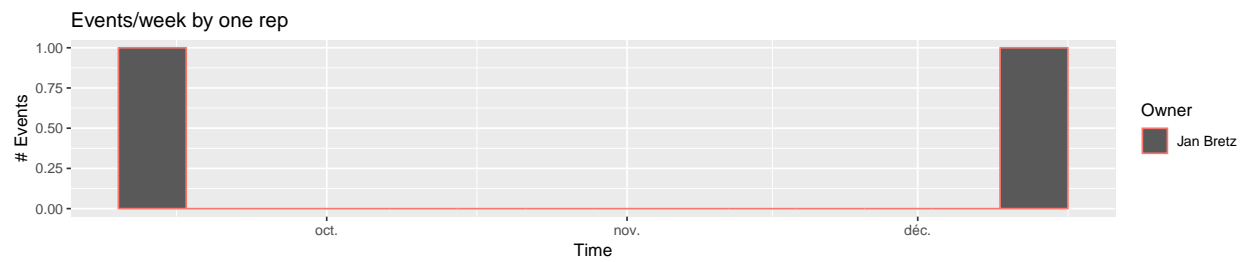
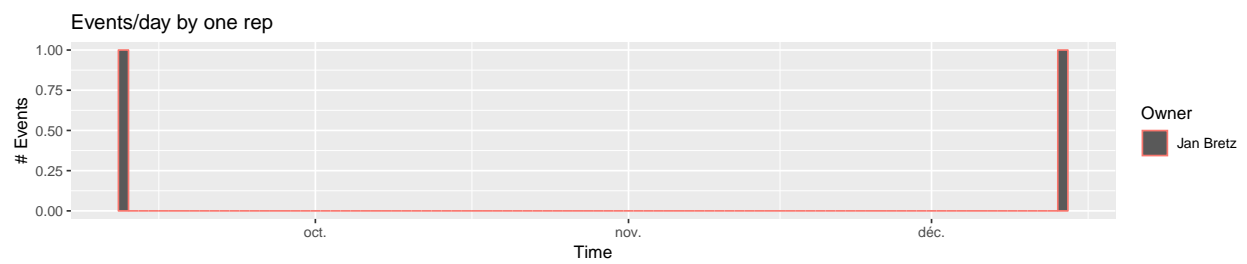
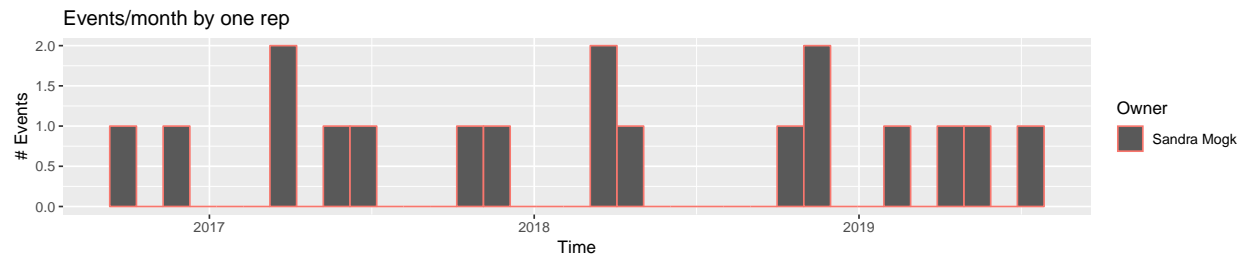
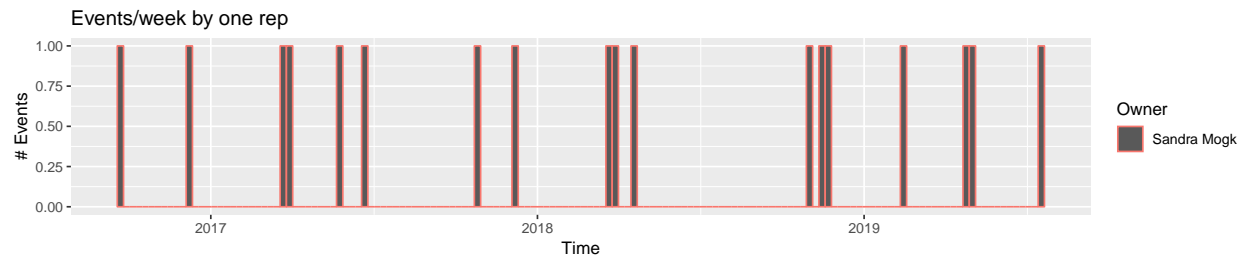
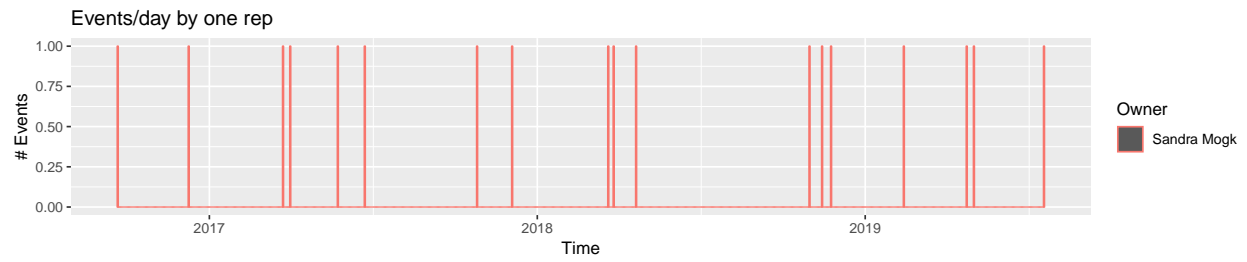


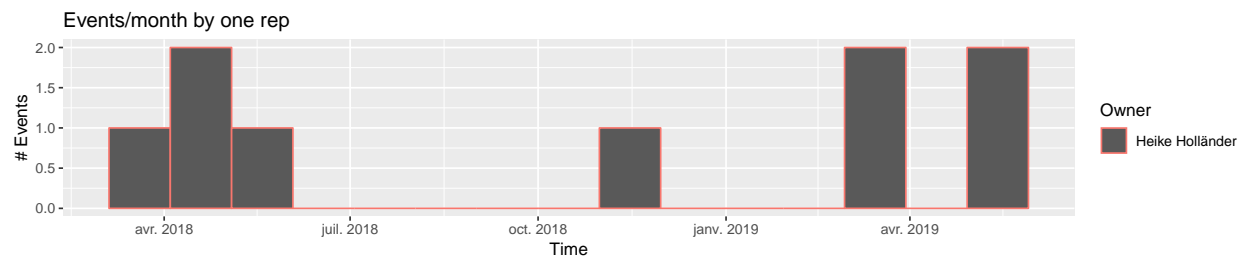
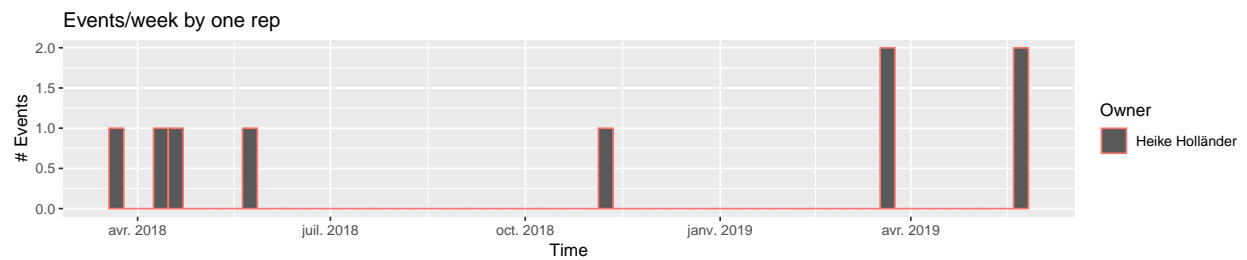
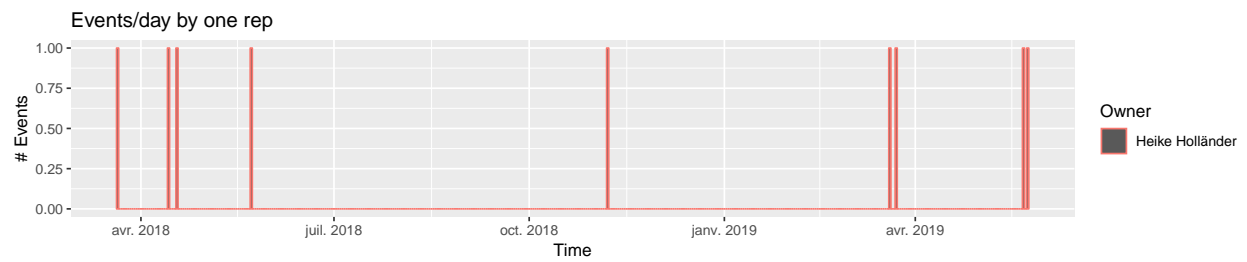
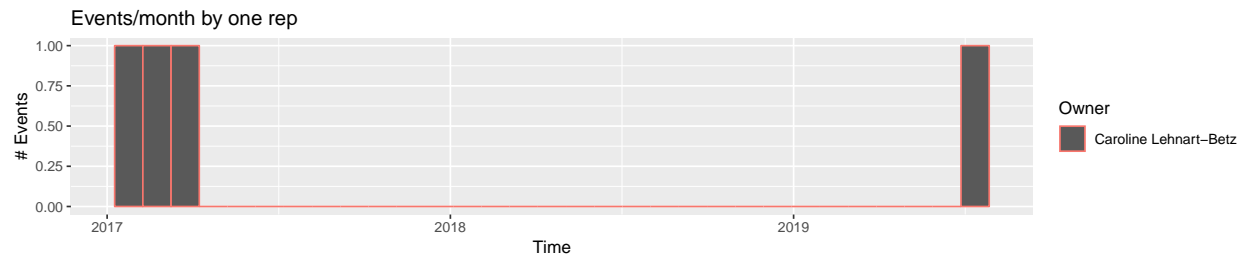
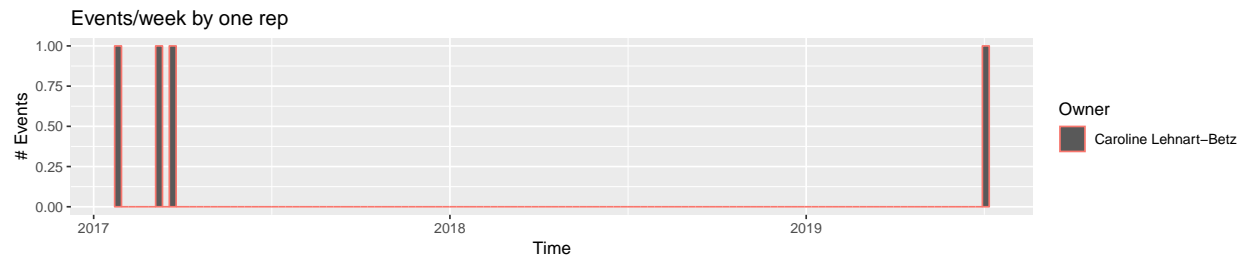
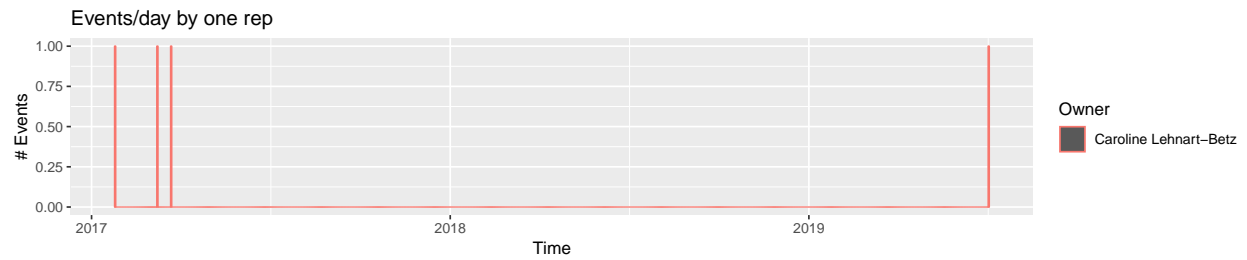


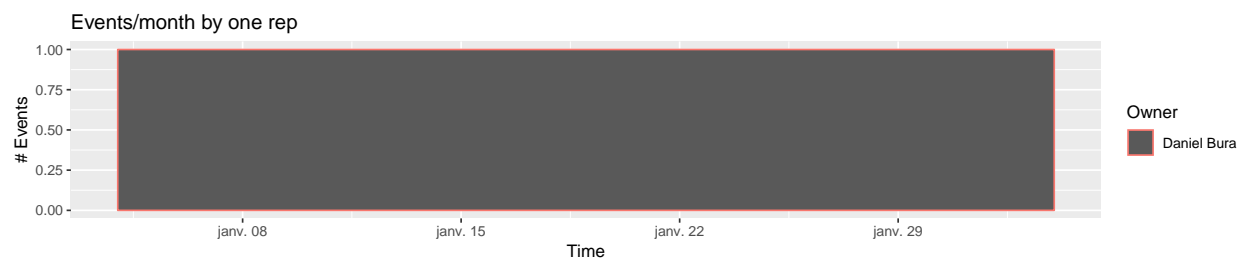
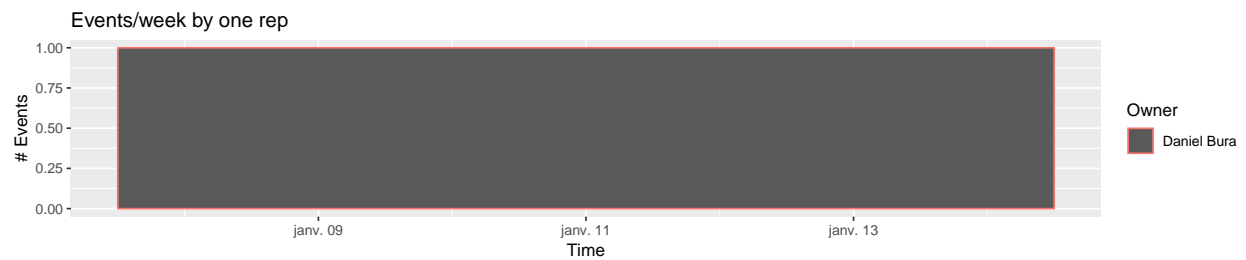
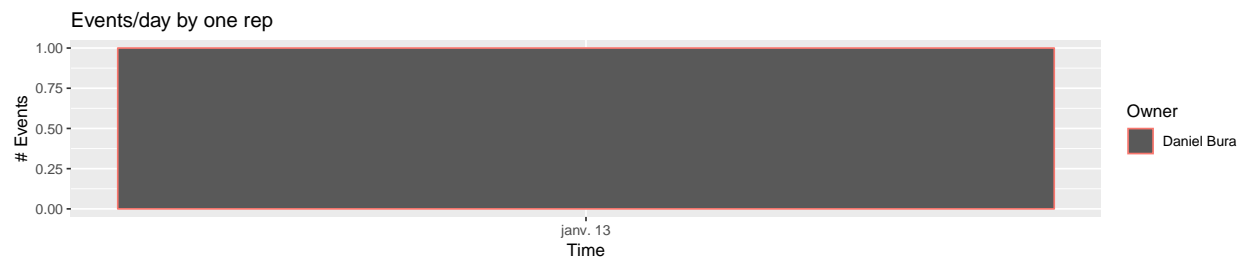
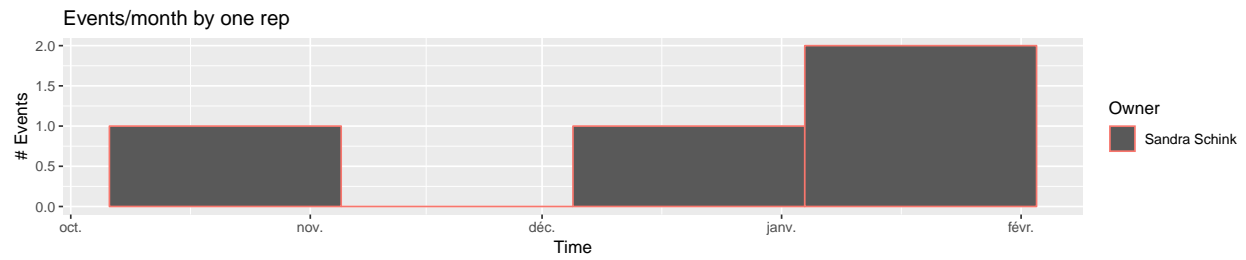
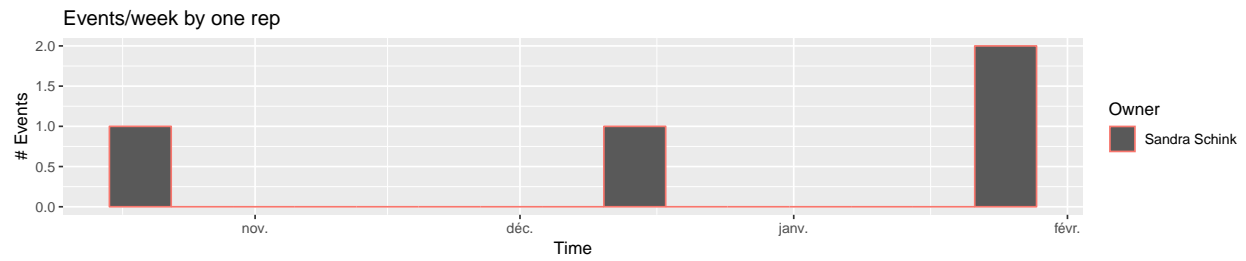
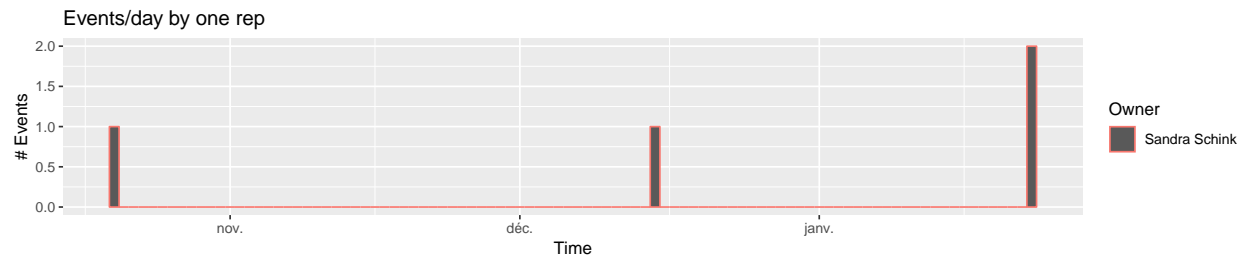


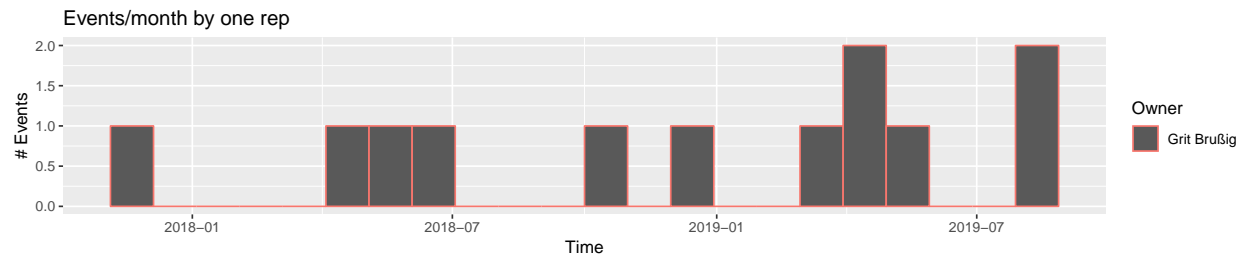
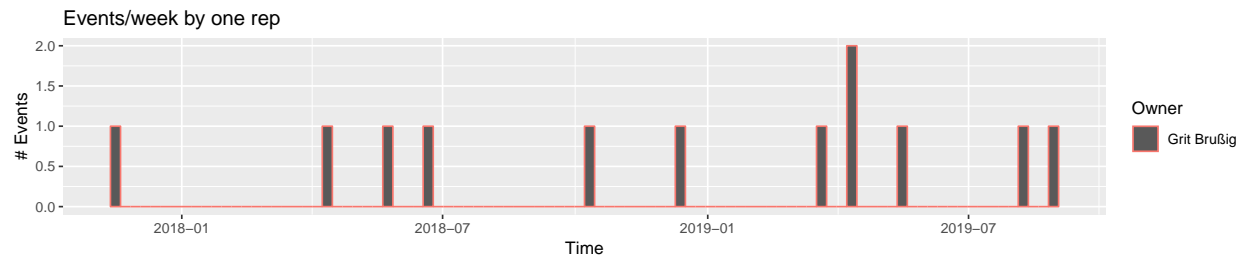
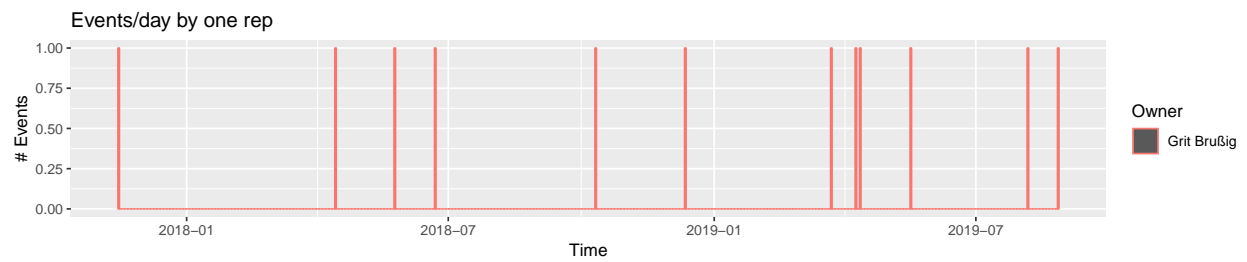
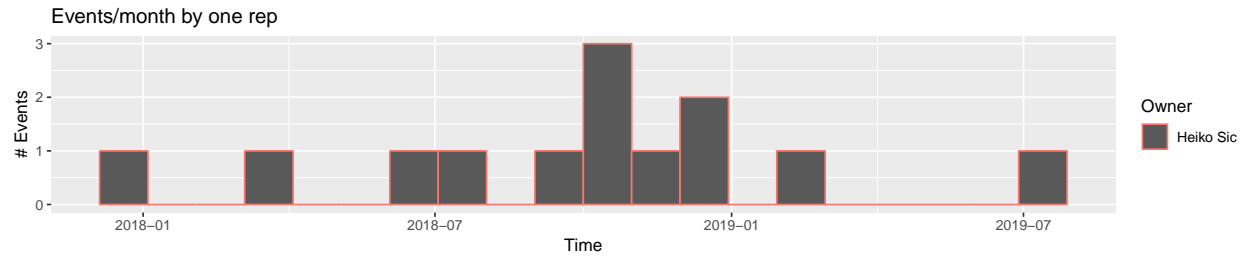
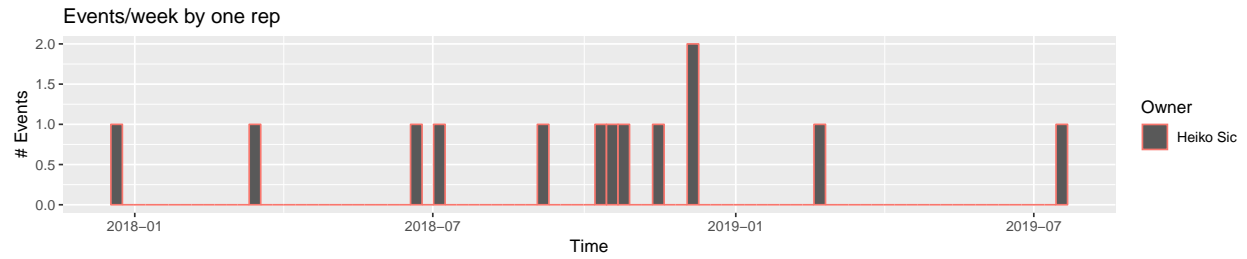
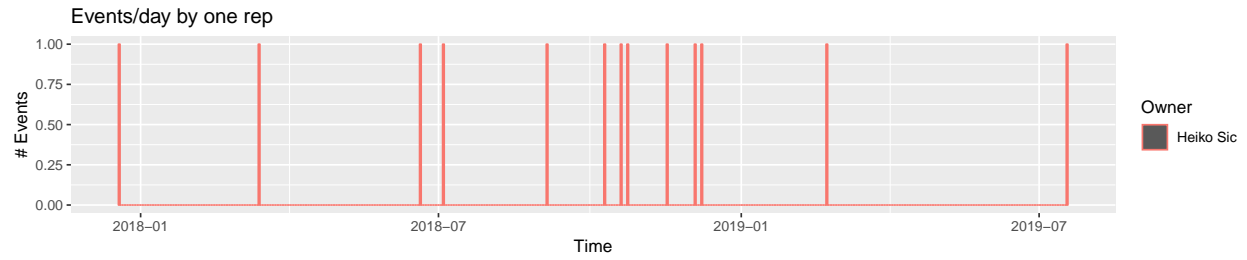


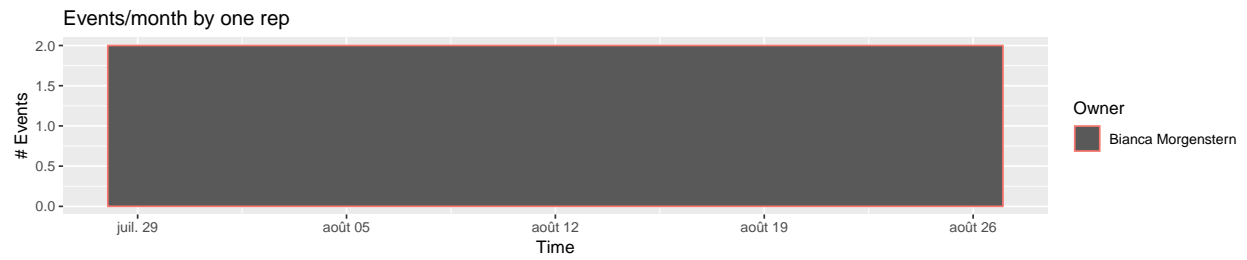
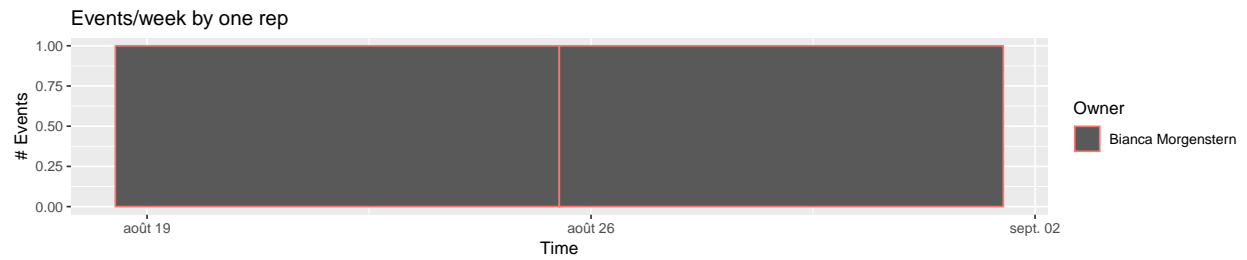
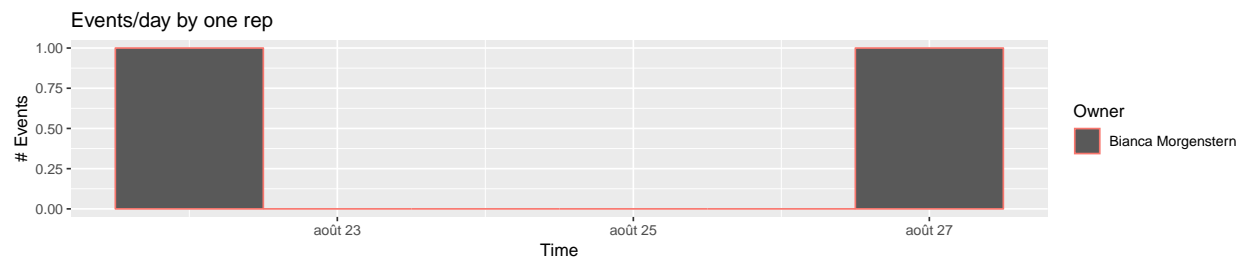
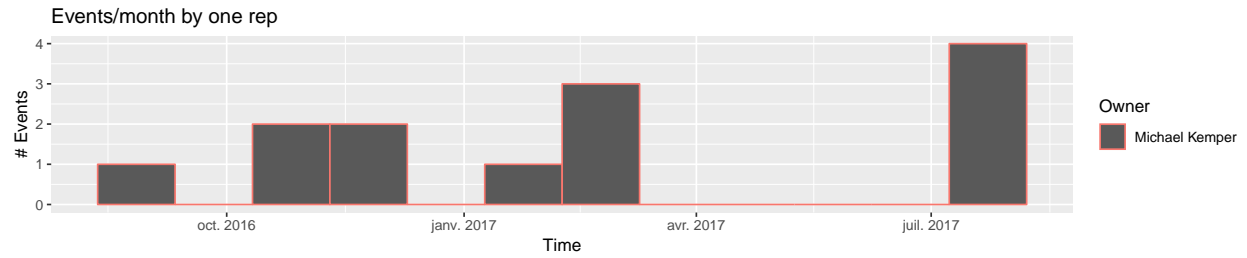
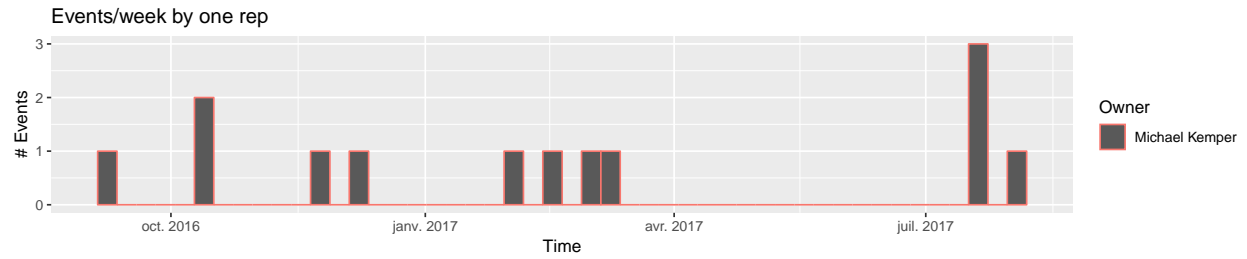
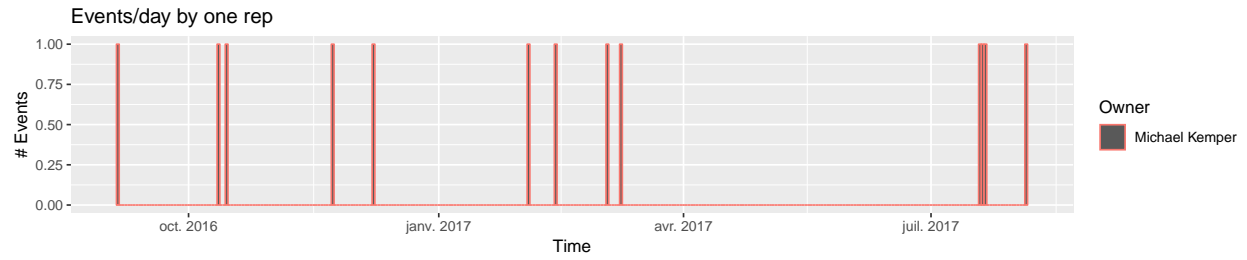


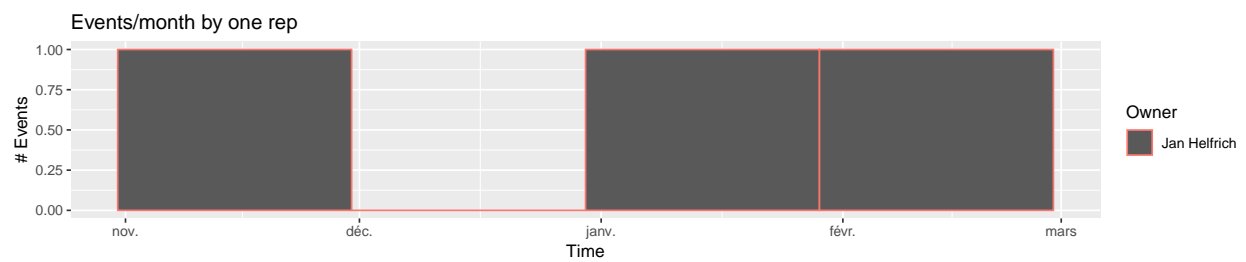
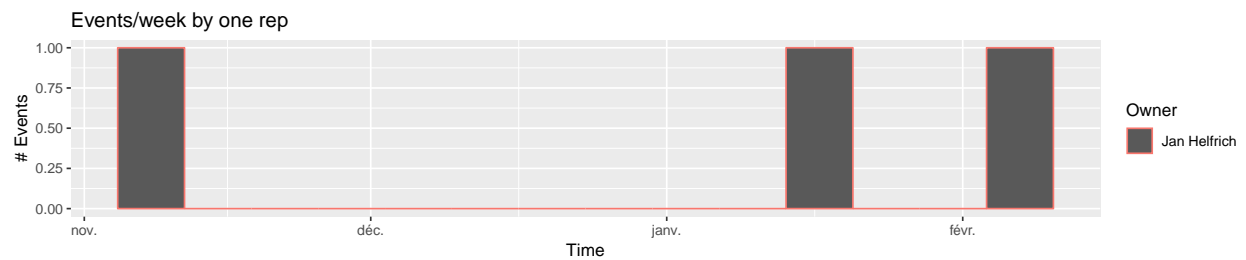
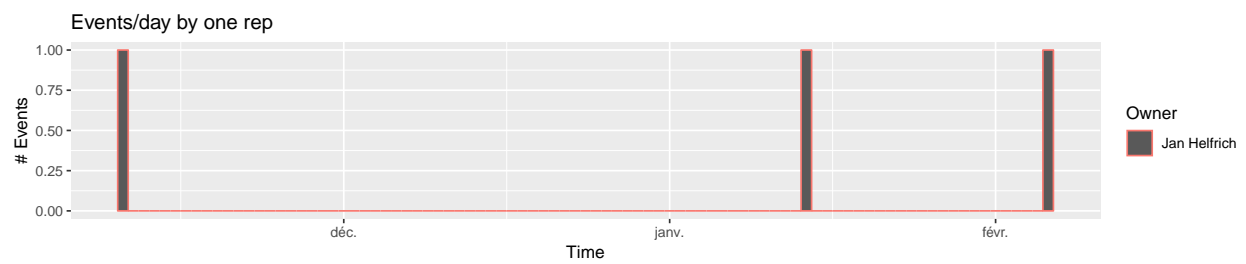
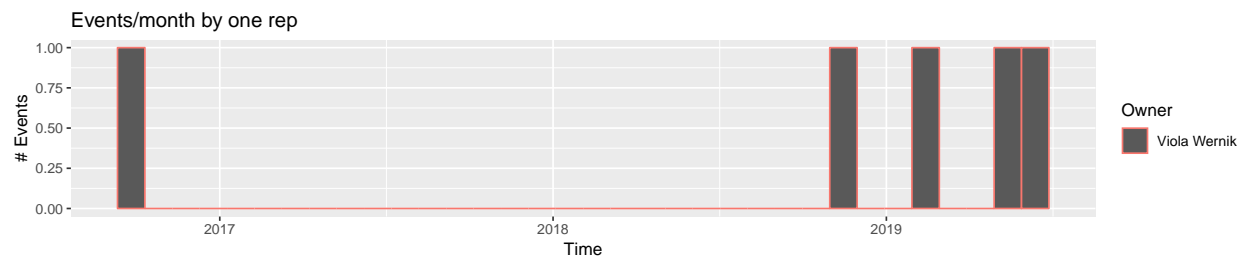
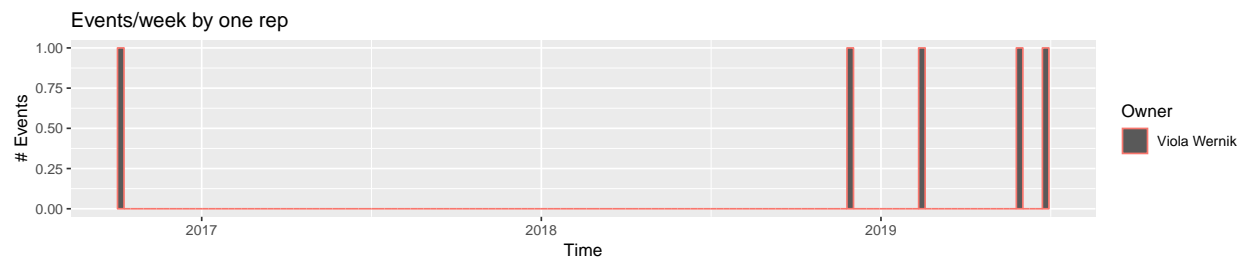
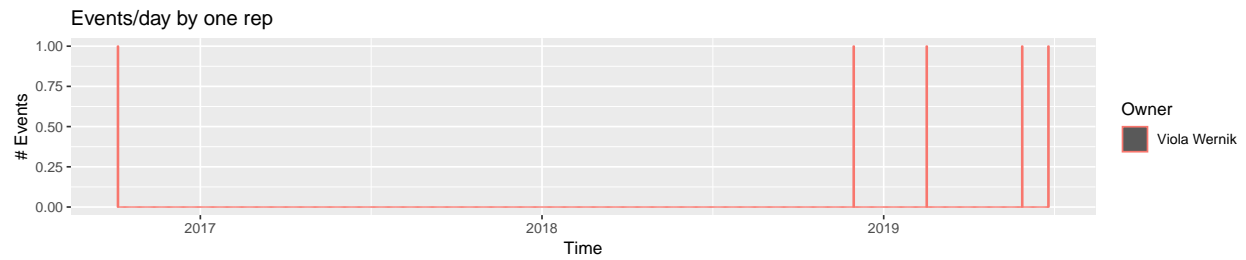


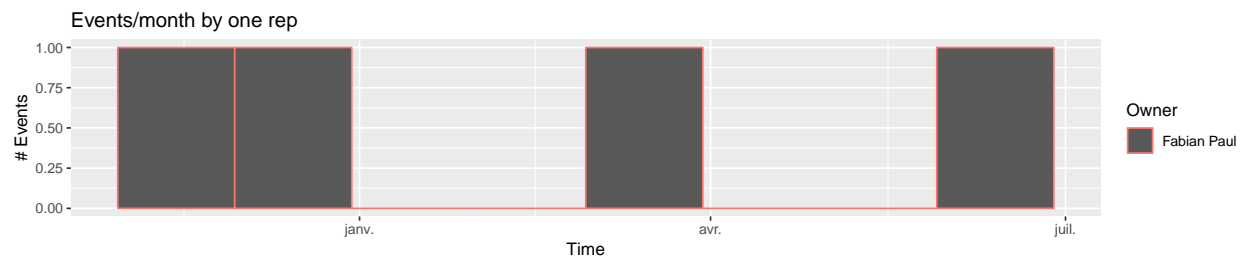
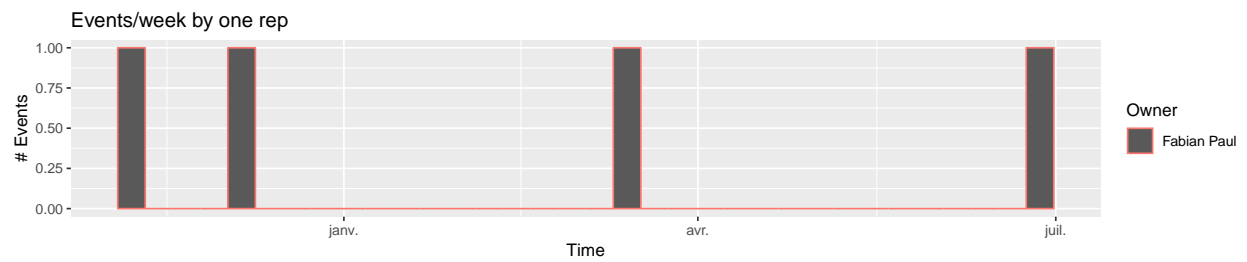
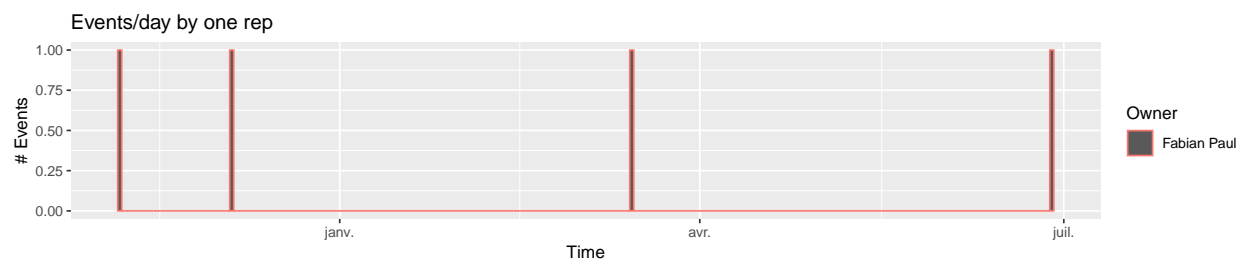
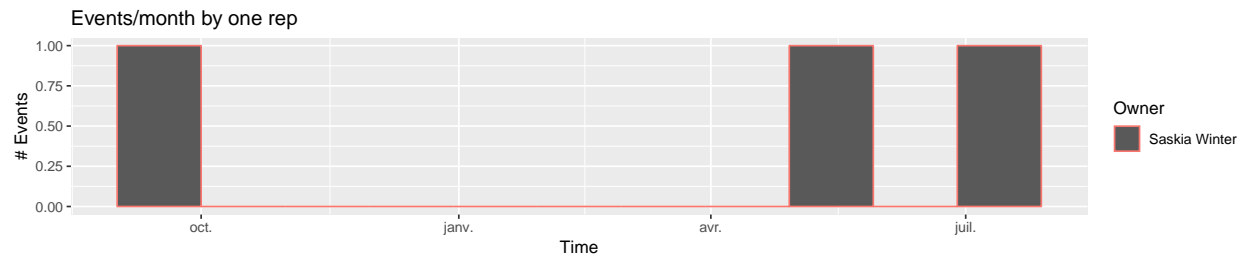
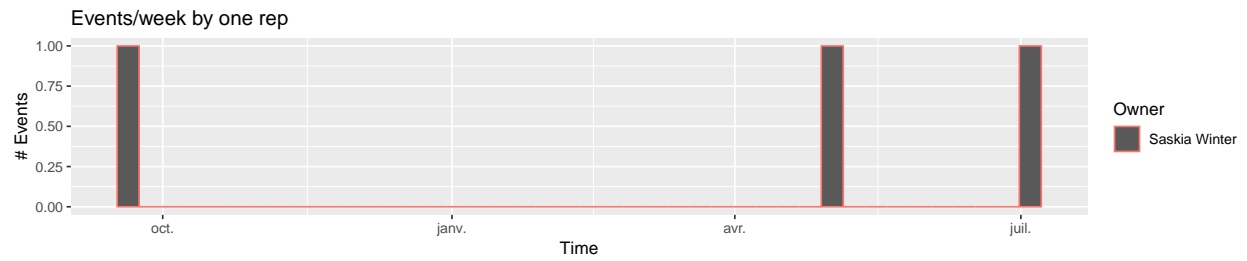
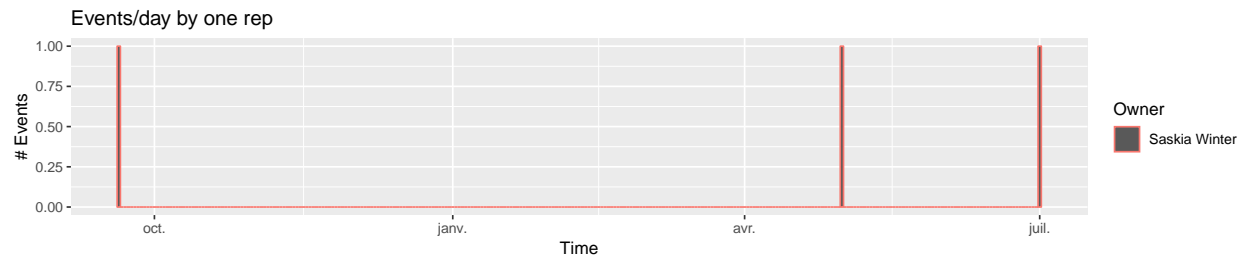


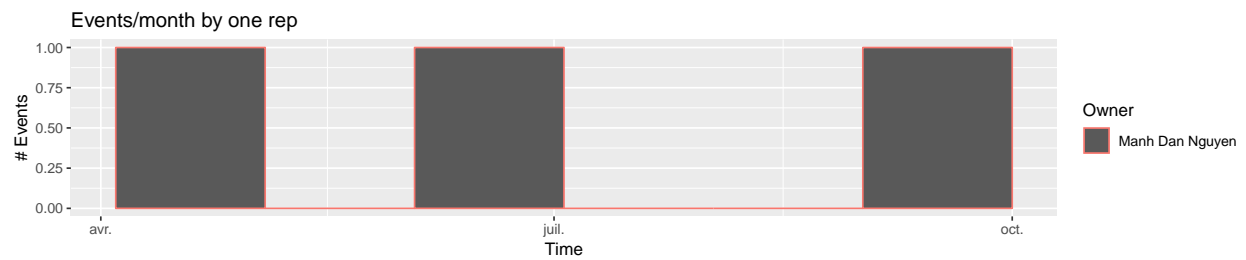
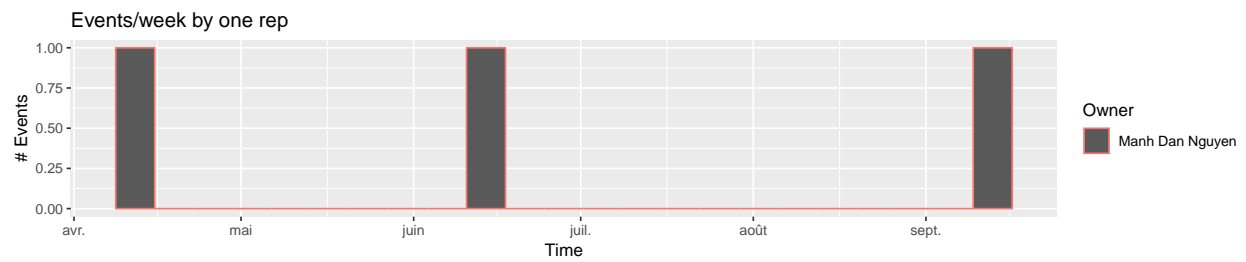
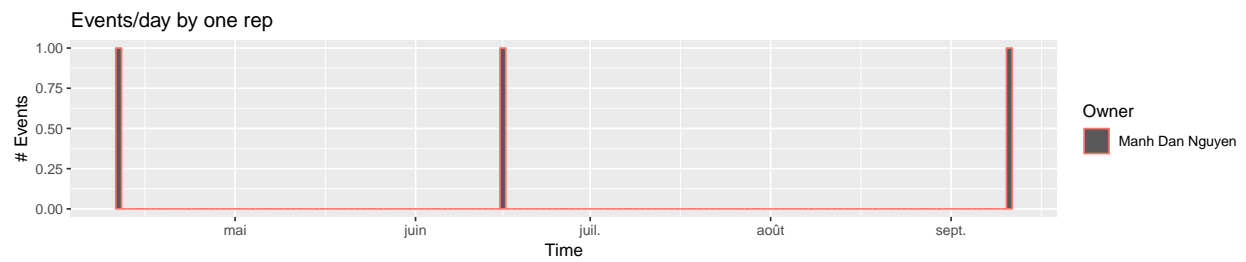
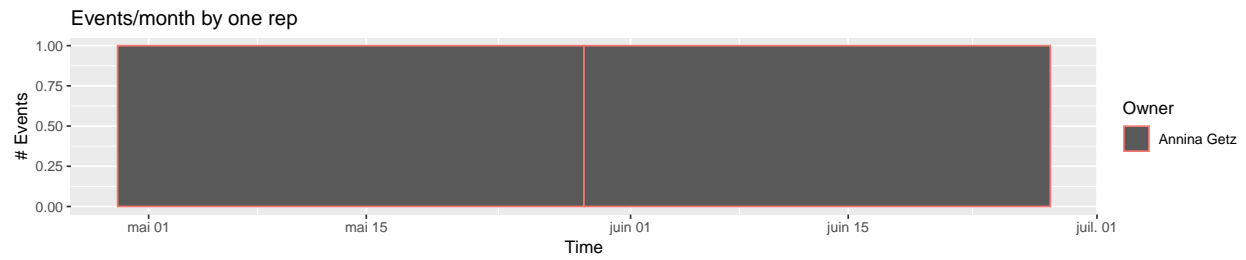
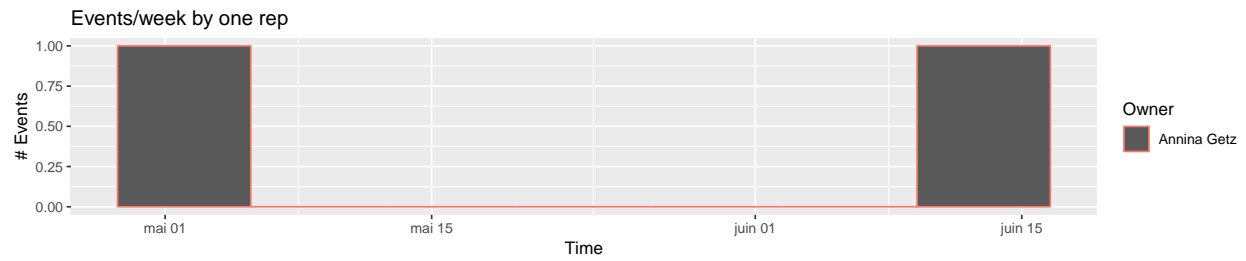
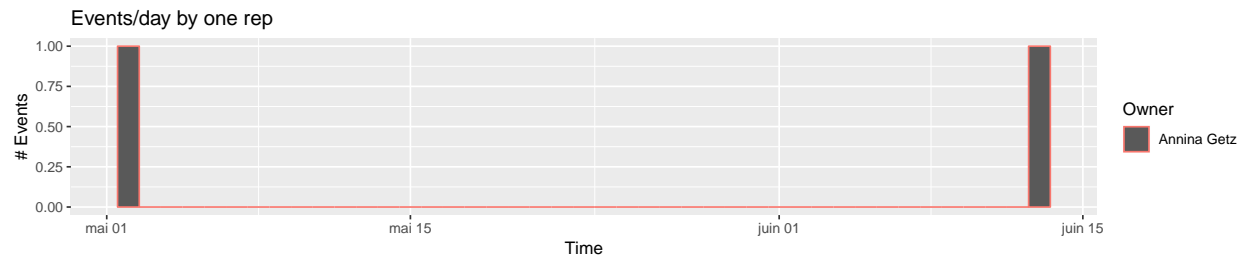


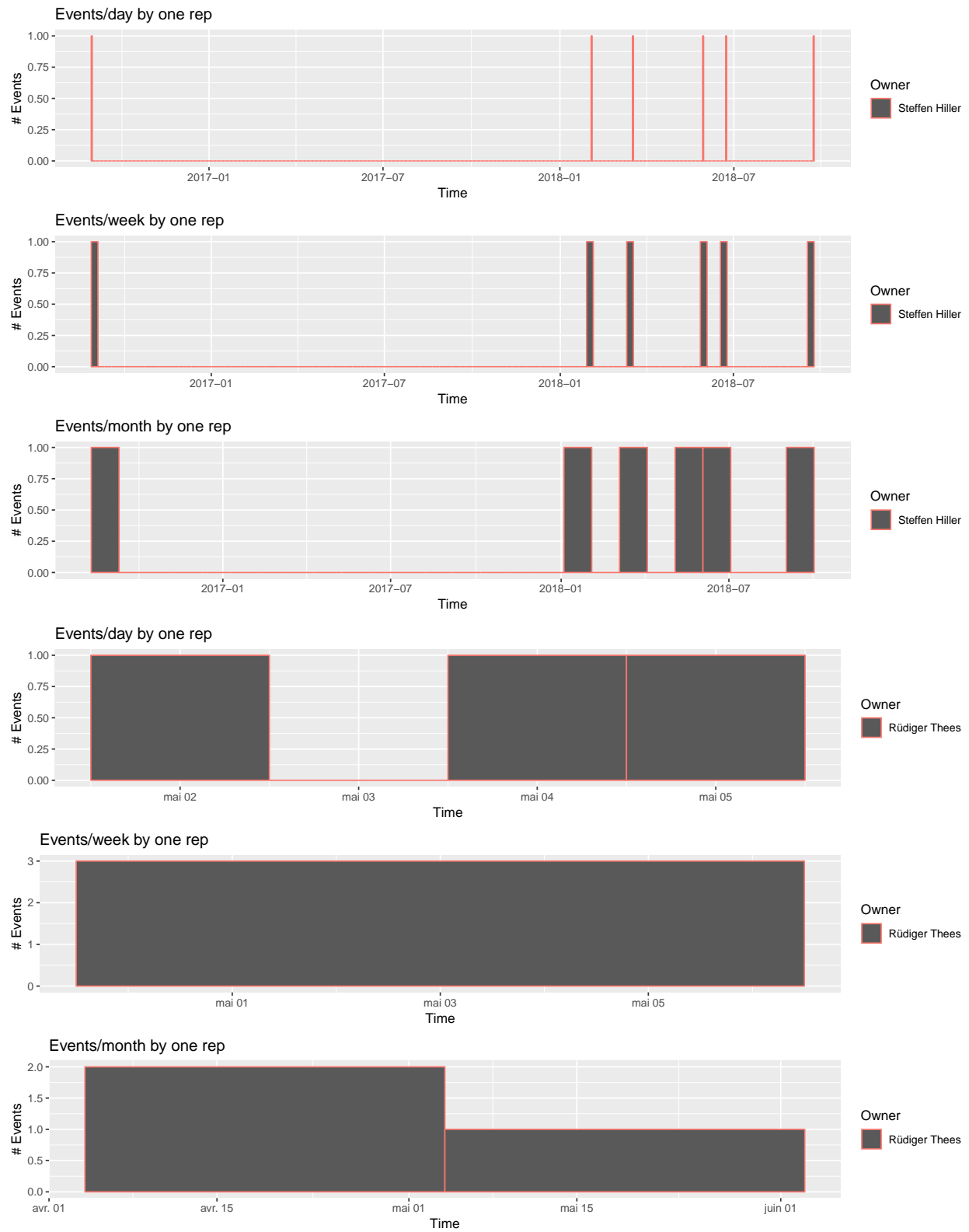


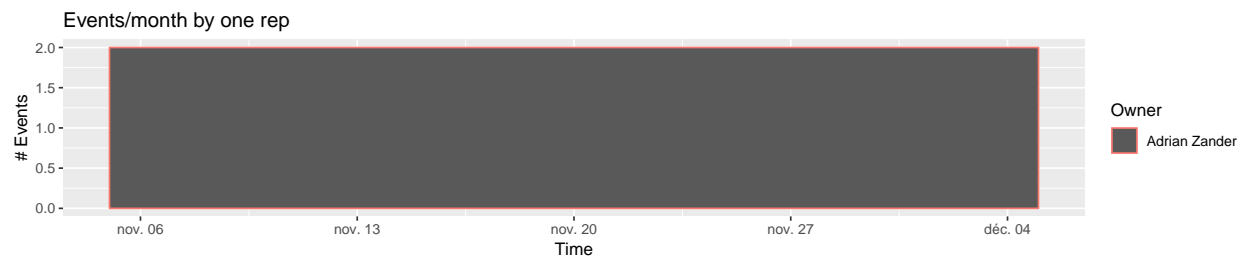
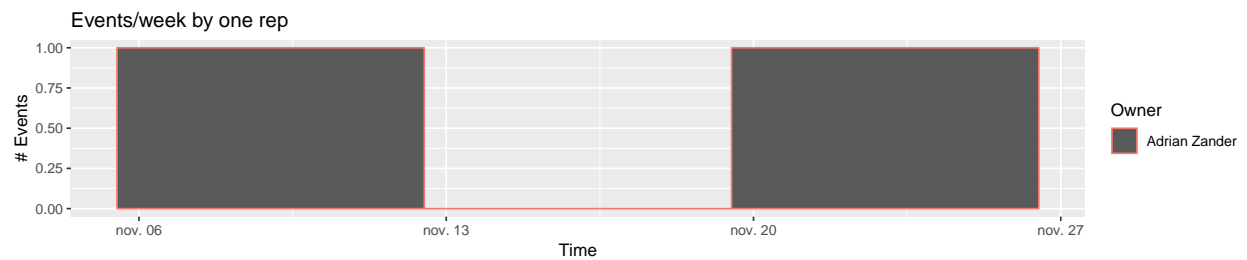
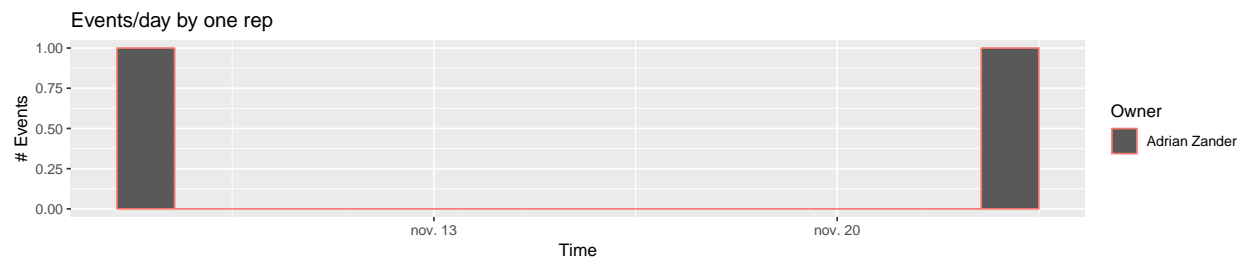
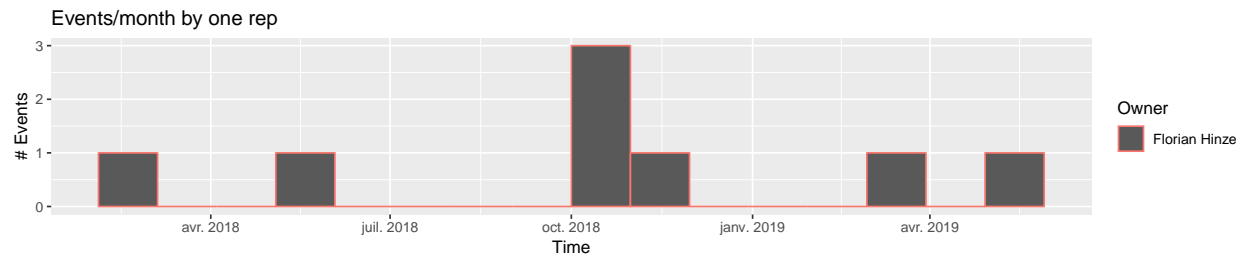
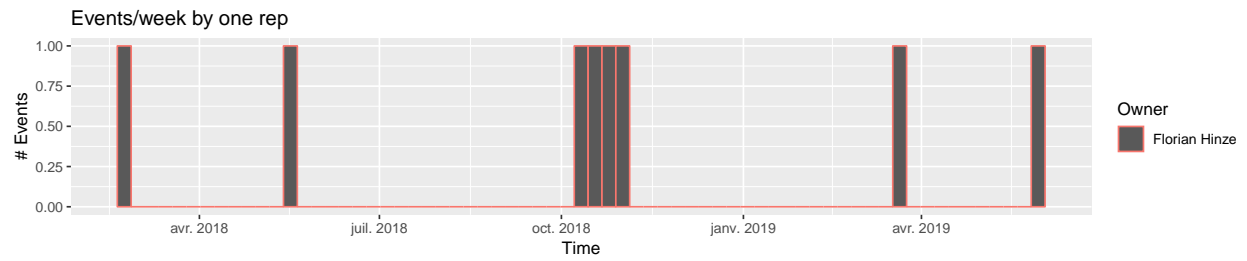
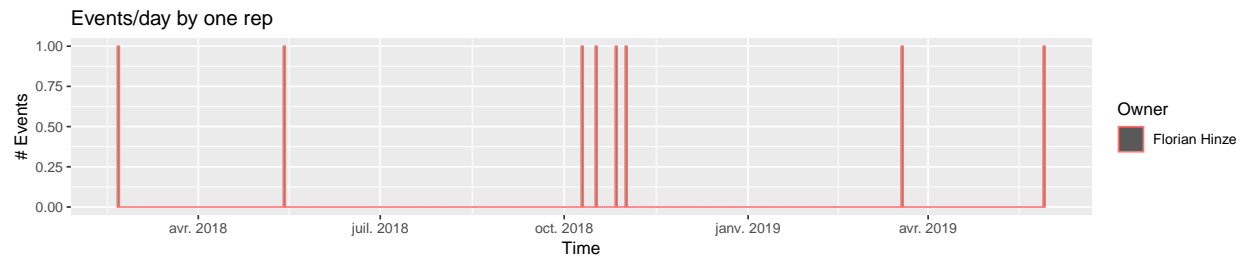


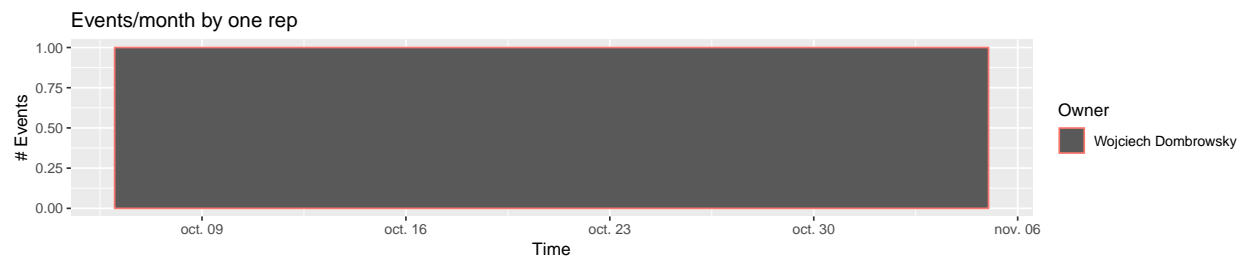
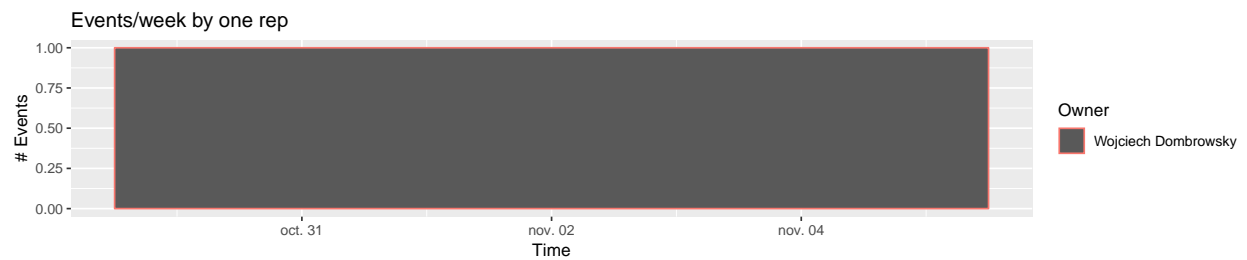
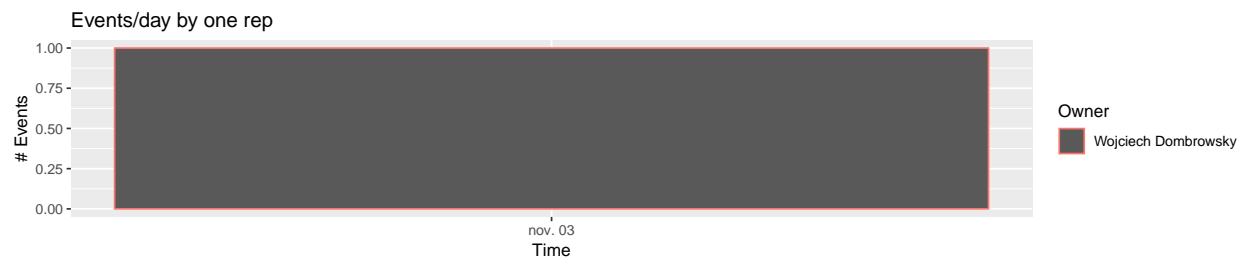
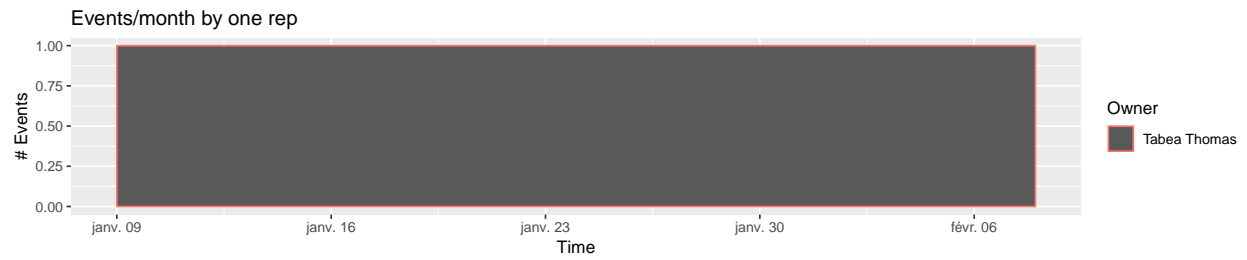
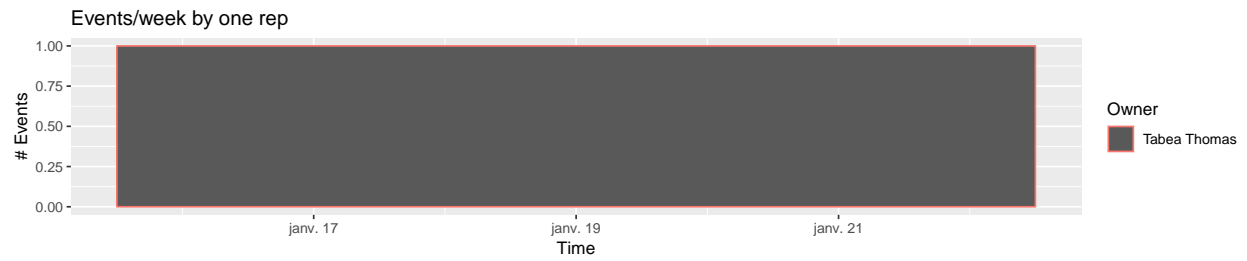
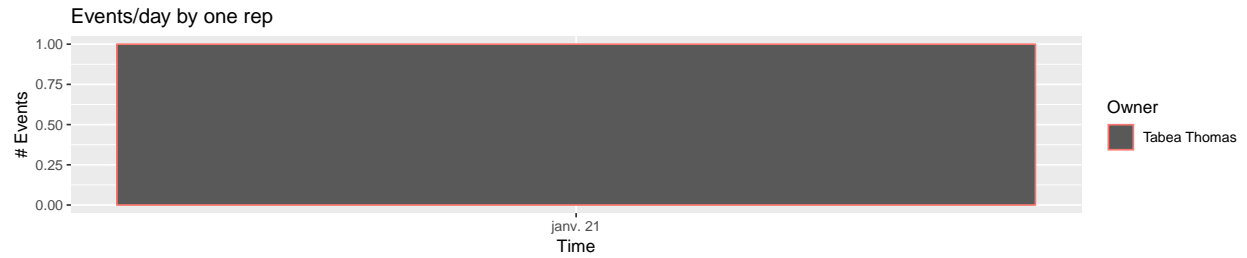


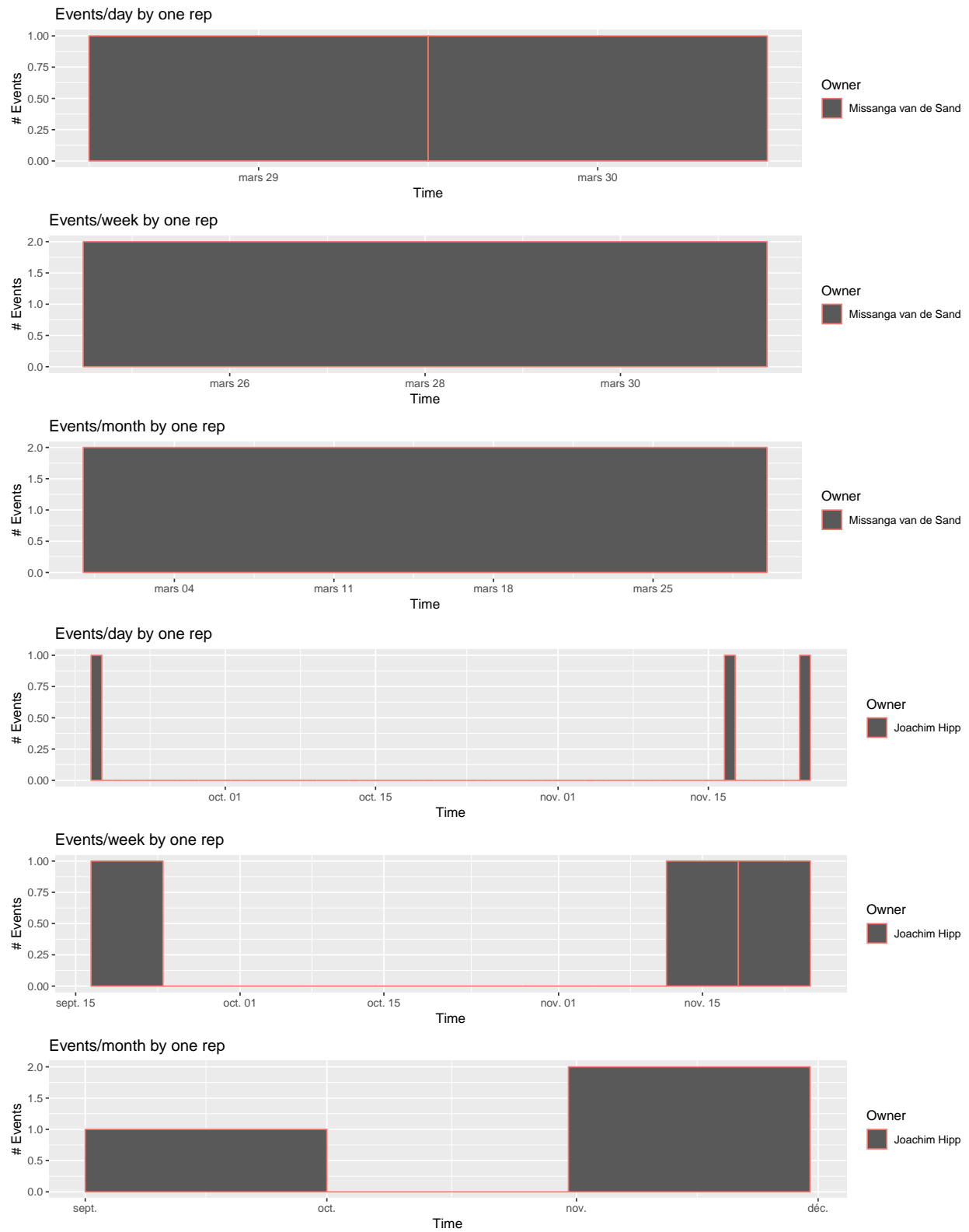








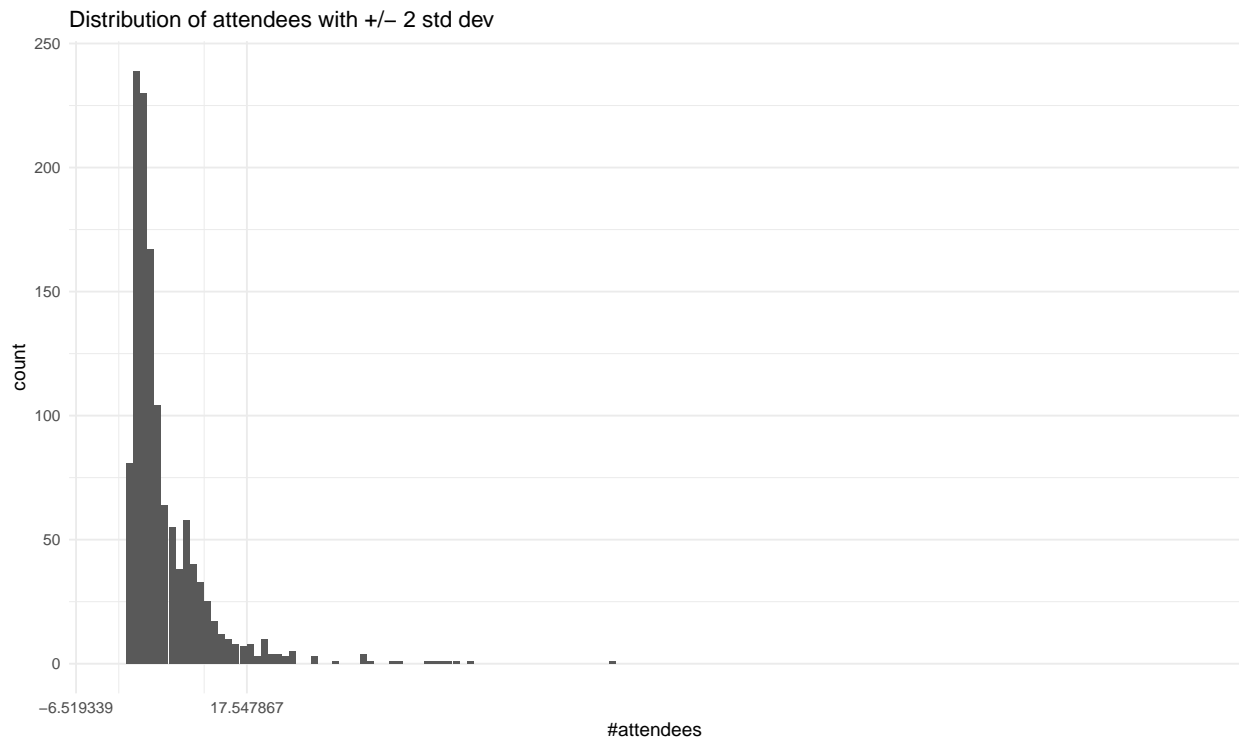




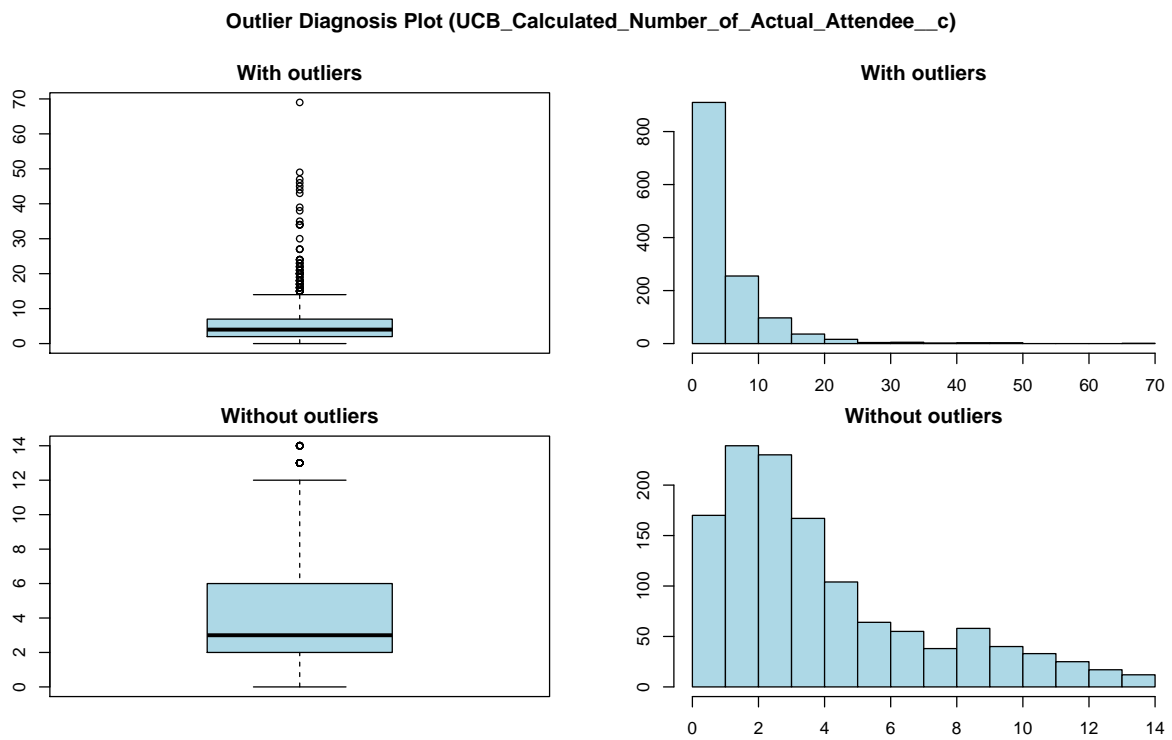
Distribution of attendees each event

[1] "mean: 5.51426426426426 std dev: 6.01680148079842"

pdf
2

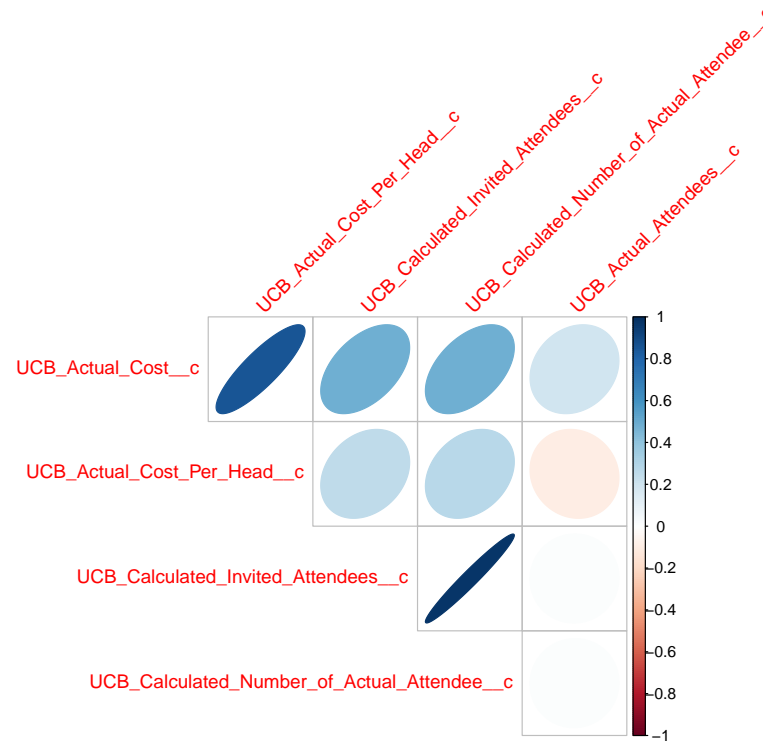


pdf
2



Bivariate data EDA. Calculation of correlation coefficient using `correlate()`

Visualization of the correlation matrix using `plot_correlate()` `plot_correlate()` visualizes the correlation matrix.



5. DQA Dimension5/Rule5: understanding consistency of data from more than one source

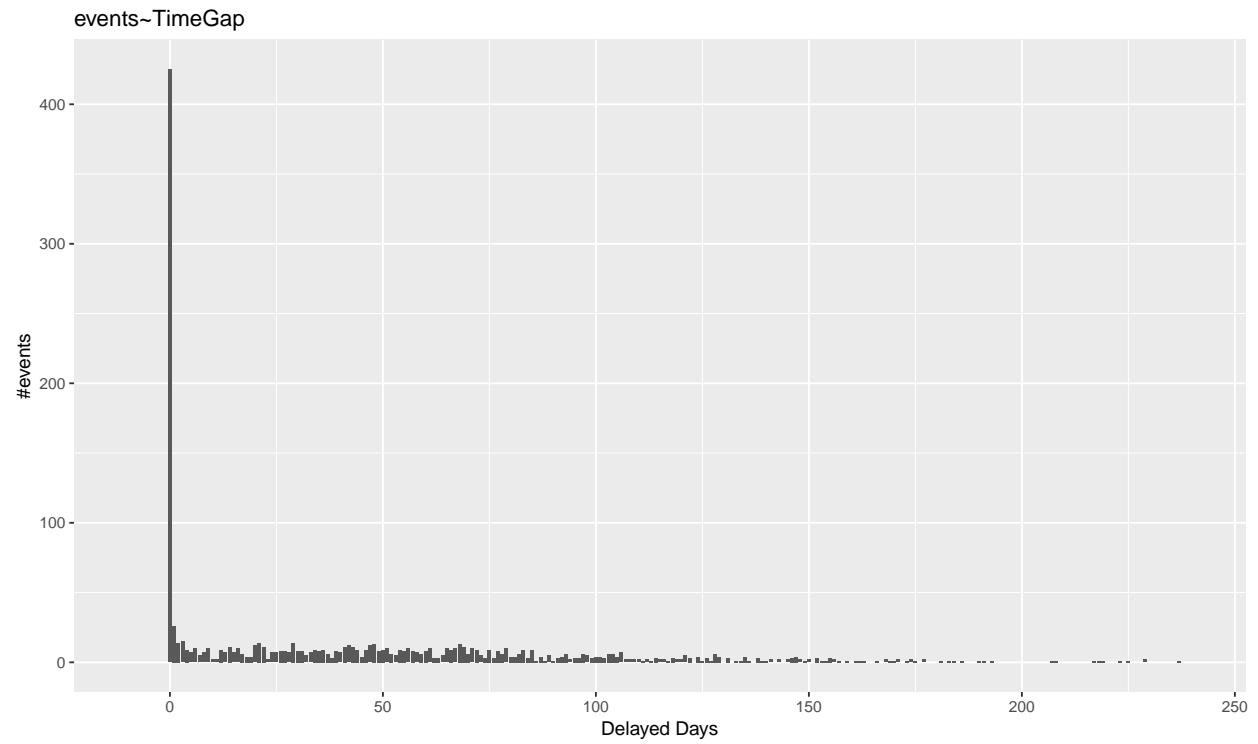
Need other data sources for comparsion

6. DQA Dimension6/Rule6: undestanding timeline, trend & seasonality

The degree to which data represent reality from the required point in time Time difference: $\text{Ind_6} = \max[(\text{threshold} - \text{xx days}) / \text{threshold}, 0]$

pdf

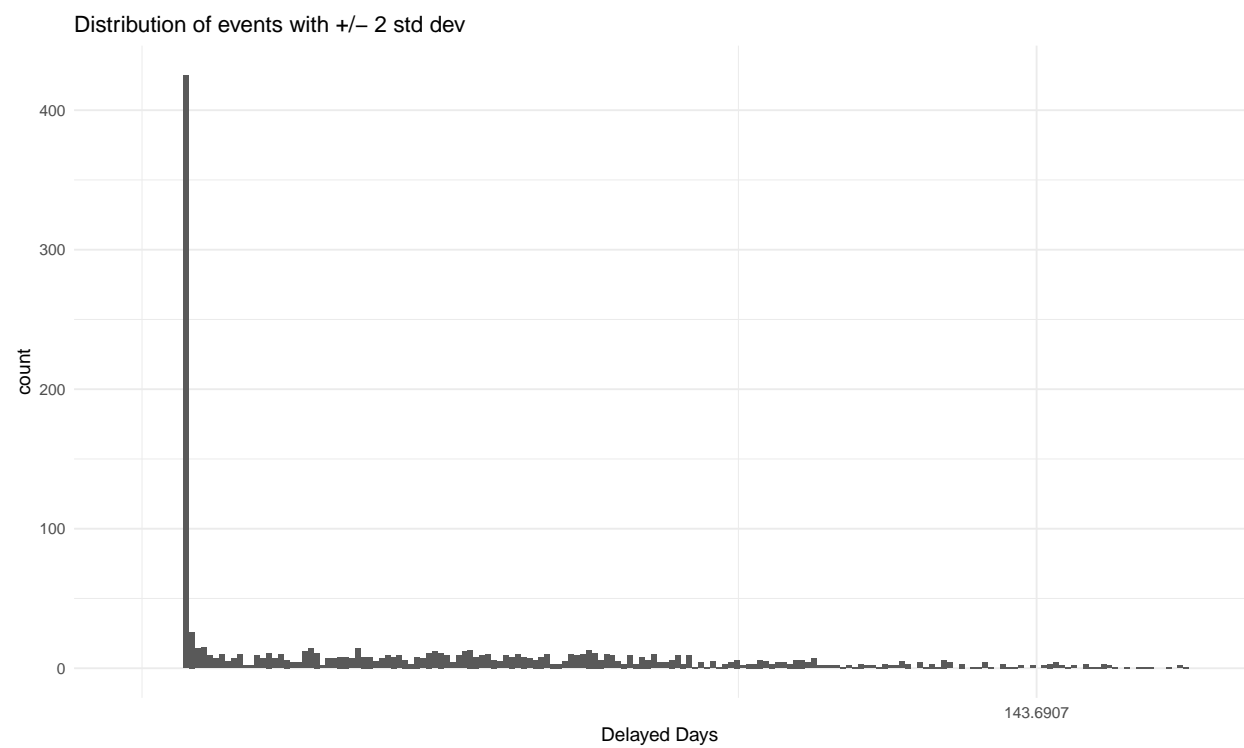
2



[1] "mean: 42.8843843843844 std dev: 50.4031544932102"

pdf

2



pdf

Outlier Diagnosis Plot (TimeGap)

