

Processing, cleaning and saving NZ GREEN Grid project 1 minute electricity power data

Ben Anderson (b.anderson@soton.ac.uk, @dataknut)

Last run at: 2018-05-18 16:40:30

Contents

1	Status	2
2	Citation	2
3	Introduction	3
3.1	Purpose	3
3.2	Requirements:	3
3.3	History	3
3.4	Support	3
4	Obtain listing of files	3
4.1	Date format checks	5
4.2	Data file quality checks	7
5	Load data files	9
6	Data quality analysis	10
6.1	Circuit label checks	10
6.2	Observations	15
7	Runtime	18
8	R environment	18
	References	19

1 Status

Full run using all data from /Volumes/hum-csafe/Research Projects/GREEN Grid/_RAW DATA/GridSpyData/

2 Citation

If you wish to use any of the material from this report please cite as:

- Anderson, B. (2018) Processing, cleaning and saving NZ GREEN Grid project 1 minute electricity power data, University of Otago: Dunedin, NZ.

3 Introduction

Report circulation:

- Restricted to: NZ GREEN Grid project partners and contractors.

3.1 Purpose

This report is intended to:

- load and clean the project electricity power data (Grid Spy)
- save the cleaned data out as a single file per household
- produce summary data quality statistics

The resulting cleaned data has *no* identifying information such as names, addresses, email addresses, telephone numbers and is therefore safe to share across all partners.

The data contains a unique household id which can be used to link it to the NZ GREEN Grid time use diaries and dwelling/appliance surveys. With some additional non-disclosure checks it should also be safe to archive all of these linkable datasets for re-use via the UK reshare service.

3.2 Requirements:

- grid spy 1 minute data downloads

3.3 History

Generally tracked via our git.soton repo:

- history
- issues

3.4 Support

This work was supported by:

- The University of Otago
- The New Zealand Ministry of Business, Innovation and Employment (MBIE)
- SPATIALEC - a Marie Skłodowska-Curie Global Fellowship based at the University of Otago's Centre for Sustainability (2017-2019) & the University of Southampton's Sustainable Energy Research Group (2019-2022).

This work is (c) 2018 the University of Southampton.

We do not 'support' the code but if you have a problem check the issues on our repo and if it doesn't already exist, open one. We might be able to fix it :-)

4 Obtain listing of files

In this section we generate a listing of all 1 minute data files that we have received. If we are running over the complete dataset then we will be using data from:

- /hum-csafe/Research Projects/GREEN Grid/_RAW DATA/GridSpyData/

In this run we are using data from:

- /Volumes/hum-csafe/Research Projects/GREEN Grid/_RAW DATA/GridSpyData/

If these do not match then this may be a test run.

```
## Loading required package: data.table
## Loading required package: lubridate
##
## Attaching package: 'lubridate'
## The following objects are masked from 'package:data.table':
##
##     hour, isoweek, mday, minute, month, quarter, second, wday,
##     week, yday, year
## The following object is masked from 'package:base':
##
##     date
## Loading required package: readr
## Loading required package: dplyr
##
## Attaching package: 'dplyr'
## The following objects are masked from 'package:lubridate':
##
##     intersect, setdiff, union
## The following objects are masked from 'package:data.table':
##
##     between, first, last
## The following object is masked from 'package:ggplot2':
##
##     vars
## The following objects are masked from 'package:stats':
##
##     filter, lag
## The following objects are masked from 'package:base':
##
##     intersect, setdiff, setequal, union
## Loading required package: progress
## [1] "Looking for 1 minute data using pattern = *at1.csv$ in /Volumes/hum-csafe/Research Projects/GREEN Grid/_RAW DATA/GridSpyData/"
## [1] "Looking for data using pattern = *at1.csv$ in /Volumes/hum-csafe/Research Projects/GREEN Grid/_RAW DATA/GridSpyData/"
## [1] "Found 21,836 files"
## [1] "Processing file list and getting file meta-data (please be patient)"
## [1] "All files checked"
## [1] "Saving 1 minute data files interim metadata to /Volumes/hum-csafe/Research Projects/GREEN Grid/_RAW DATA/GridSpyData/"
## [1] "Done"
## [1] "Overall we have 21836 files from 44 households."
```

Overall we have 21,836 files from 44 households. Of the 21,836, 12,722 (58.26%) were *not* loaded/checked as their file sizes indicated that they contained no data.

4.1 Date format checks

We now need to check how many of the loaded files have an ambiguous or default date - these could introduce errors.

Table 1: Number of files and min/max date (as char) with given date column names by inferred date format

dateColName	dateFormat	nFiles	minDate	maxDate
date NZ	dmy - definite	1	27/03/2015	27/03/2015
date NZ	mdy - definite	2	5/26/2016	5/26/2016
date NZ	ymd - default (but day/month value <= 12)	12	2014-01-06	2016-06-07
date NZ	ymd - definite	67	2014-05-24	2016-07-13
date UTC	ambiguous	28	11-10-16	27/07/14
date UTC	ymd - default (but day/month value <= 12)	3607	2014-11-03	2018-05-12
date UTC	ymd - definite	5397	2015-03-26	2018-05-15
unknown - do not load (fsize = 2751)	NA	1812	NA	NA
unknown - do not load (fsize = 43)	NA	10910	NA	NA

Results to note:

- There are 28 ambiguous files
- The non-loaded files only have 2 distinct file sizes, confirming that they are unlikely to contain useful data.

We now inspect the ambiguous and (some of) the default files.

To help with data cleaning the following table lists files that have ambiguous dates.

```
# list ambiguous files
aList <- fListCompleteDT[dateFormat == "ambiguous",
  .(file, dateColName, dateExample, dateFormat)]

cap <- paste0("All ", nrow(aList),
  " files with an ambiguous dateFormat")

knitr::kable(caption = cap, aList)
```

Table 2: All 28 files with an ambiguous dateFormat

file	dateColName	dateExample	dateFormat
rf_06/15Jul2014-25May2016at1.csv	date UTC	14/07/14	ambiguous
rf_07/15Jul2014-25May2016at1.csv	date UTC	14/07/14	ambiguous
rf_08/15Jul2014-25May2016at1.csv	date UTC	14/07/14	ambiguous
rf_10/15Jul2014-25May2016at1.csv	date UTC	14/07/14	ambiguous
rf_11/15Jul2014-25May2016at1.csv	date UTC	14/07/14	ambiguous
rf_13/15Jul2014-25May2016at1.csv	date UTC	14/07/14	ambiguous
rf_19/15Jul2014-25May2016at1.csv	date UTC	14/07/14	ambiguous
rf_21/15Jul2014-25May2016at1.csv	date UTC	14/07/14	ambiguous
rf_22/15Jul2014-25May2016at1.csv	date UTC	14/07/14	ambiguous
rf_23/15Jul2014-25May2016at1.csv	date UTC	14/07/14	ambiguous
rf_24/15Jul2014-25May2016at1.csv	date UTC	27/07/14	ambiguous
rf_25/12Oct2016-20Nov2017at1.csv	date UTC	11-10-16	ambiguous
rf_26/15Jul2014-25May2016at1.csv	date UTC	14/07/14	ambiguous

file	dateColName	dateExample	dateFormat
rf_27/15Jul2014-25May2016at1.csv	date UTC	27/07/14	ambiguous
rf_29/24Mar2015-25May2016at1.csv	date UTC	25/03/15	ambiguous
rf_30/15Feb2016-25May2016at1.csv	date UTC	14/02/16	ambiguous
rf_30/24Mar2015-25May2016at1.csv	date UTC	27/03/15	ambiguous
rf_31/24Mar2015-25May2016at1.csv	date UTC	25/03/15	ambiguous
rf_34/18Jan2016-25May2016at1.csv	date UTC	17/01/16	ambiguous
rf_34/20Jul2015-25May2016at1.csv	date UTC	19/07/15	ambiguous
rf_34/24Mar2015-25May2016at1.csv	date UTC	26/03/15	ambiguous
rf_35/24Mar2015-25May2016at1.csv	date UTC	23/03/15	ambiguous
rf_39/24Mar2015-25May2016at1.csv	date UTC	27/03/15	ambiguous
rf_43/24Mar2015-25May2016at1.csv	date UTC	26/03/15	ambiguous
rf_43/27Mar2015-18Oct2015at1.csv	date UTC	26/03/15	ambiguous
rf_44/24Mar2015-25May2016at1.csv	date UTC	24/03/15	ambiguous
rf_46/12Oct2016-20Nov2017at1.csv	date UTC	11-10-16	ambiguous
rf_47/24Mar2015-25May2016at1.csv	date UTC	24/03/15	ambiguous

Check against file names to see what is reasonable and then fix them.

```
# Setting to dmy seems OK
fListCompleteDT <- fListCompleteDT[dateFormat == "ambiguous",
                                     dateFormat := "dmy - inferred"]

paste0("Fixed ", nrow(aList), " files with an ambiguous dateFormat")
```

```
## [1] "Fixed 28 files with an ambiguous dateFormat"
```

The following table lists up to 10 of the ‘date NZ’ files which are set by default - do they look OK to assume the default dateFormat? Compare the file names with the dateExample...

```
# list default files with NZ time
aList <- fListCompleteDT[dateColName == "date NZ" & dateFormat %like% "default",
                          .(file, fSize, dateColName, dateExample, dateFormat)]

cap <- paste0("First 10 (max) of ", nrow(aList),
              " files with dateColName = 'date NZ' and default dateFormat")

knitr::kable(caption = cap, head(aList))
```

Table 3: First 10 (max) of 12 files with dateColName = ‘date NZ’ and default dateFormat

file	fSize	dateColName	dateExample	dateFormat
rf_01/1Jan2014-24May2014at1.csv	6255737	date NZ	2014-01-06	ymd - default (but day/month value <=)
rf_02/1Jan2014-24May2014at1.csv	6131625	date NZ	2014-03-03	ymd - default (but day/month value <=)
rf_06/24May2014-24May2015at1.csv	19398444	date NZ	2014-06-09	ymd - default (but day/month value <=)
rf_10/24May2014-24May2015at1.csv	24386048	date NZ	2014-07-09	ymd - default (but day/month value <=)
rf_11/24May2014-24May2015at1.csv	23693893	date NZ	2014-07-08	ymd - default (but day/month value <=)
rf_12/24May2014-24May2015at1.csv	21191785	date NZ	2014-07-09	ymd - default (but day/month value <=)

The following table lists up to 10 of the ‘date UTC’ files which are set by default - do they look OK to assume the default dateFormat? Compare the file names with the dateExample...

```
# list default files with UTC time
aList <- fListCompleteDT[dateColName == "date UTC" & dateFormat %like% "default",
                        .(file, fSize, dateColName, dateExample, dateFormat)]

cap <- paste0("First 10 (max) of ", nrow(aList),
             " files with dateColName = 'date UTC' and default dateFormat")

knitr::kable(caption = cap, head(aList, 10))
```

Table 4: First 10 (max) of 3607 files with dateColName = ‘date UTC’ and default dateFormat

file	fSize	dateColName	dateExample	dateFormat
rf_06/10Apr2018-11Apr2018at1.csv	156944	date UTC	2018-04-09	ymd - default (but day/month value <= 1
rf_06/10Dec2017-11Dec2017at1.csv	156601	date UTC	2017-12-09	ymd - default (but day/month value <= 1
rf_06/10Feb2018-11Feb2018at1.csv	153353	date UTC	2018-02-09	ymd - default (but day/month value <= 1
rf_06/10Jan2018-11Jan2018at1.csv	153982	date UTC	2018-01-09	ymd - default (but day/month value <= 1
rf_06/10Mar2018-11Mar2018at1.csv	156471	date UTC	2018-03-09	ymd - default (but day/month value <= 1
rf_06/10May2018-11May2018at1.csv	156683	date UTC	2018-05-09	ymd - default (but day/month value <= 1
rf_06/10Nov2017-11Nov2017at1.csv	155639	date UTC	2017-11-09	ymd - default (but day/month value <= 1
rf_06/11Apr2018-12Apr2018at1.csv	157181	date UTC	2018-04-10	ymd - default (but day/month value <= 1
rf_06/11Dec2017-12Dec2017at1.csv	157814	date UTC	2017-12-10	ymd - default (but day/month value <= 1
rf_06/11Feb2018-12Feb2018at1.csv	153859	date UTC	2018-02-10	ymd - default (but day/month value <= 1

Check final date formats:

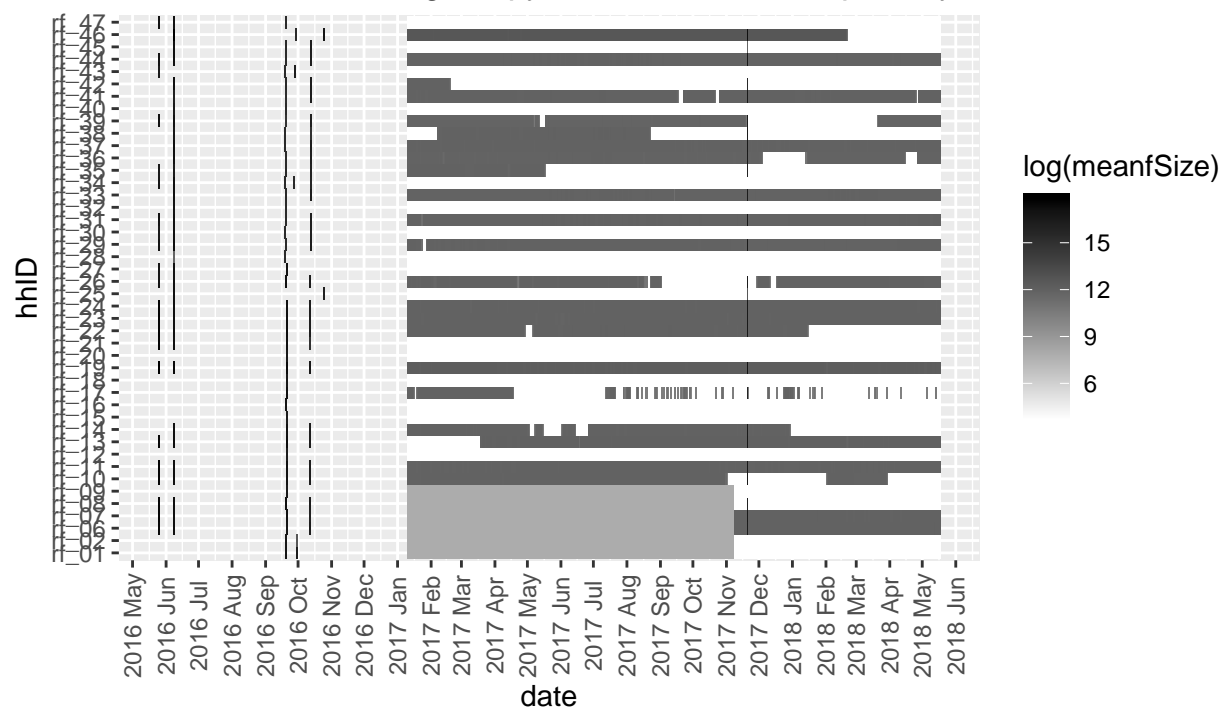
Table 5: Number of files & min/max dates (as char) with given date column names by final imputed date format

dateColName	dateFormat	nFiles	minDate	maxDate
date NZ	dmy - definite	1	27/03/2015	27/03/2015
date NZ	mdy - definite	2	5/26/2016	5/26/2016
date NZ	ymd - default (but day/month value <= 12)	12	2014-01-06	2016-06-07
date NZ	ymd - definite	67	2014-05-24	2016-07-13
date UTC	dmy - inferred	28	11-10-16	27/07/14
date UTC	ymd - default (but day/month value <= 12)	3607	2014-11-03	2018-05-12
date UTC	ymd - definite	5397	2015-03-26	2018-05-15
unknown - do not load (fsize = 2751)	NA	1812	NA	NA
unknown - do not load (fsize = 43)	NA	10910	NA	NA

4.2 Data file quality checks

The following chart shows the distribution of these files over time using their sizes. Note that white indicates the presence of small files which may not contain observations.

Mean file size of all grid spy data files received per day



ata source: /Volumes/hum-csafe/Research Projects/GREEN Grid/_RAW DATA/GridSpyData/
 Using data received up to 2018-05-18
 Log file size used as some files are full year data

Saving 6.5 x 4.5 in image

The following chart shows the same chart but only for files which we think contain data.


```
## Saving 6.5 x 4.5 in image
```

In this section we load the data files that have a file size > 3000 bytes. Things to note:

- The following table shows the number of files per household that we will load.

Table 6: Summary of household files to load

9

hhID	nFiles	meanSize	minFileDate	maxFileDate
rf_02	493	1701.856	2017-01-11	2018-05-17
rf_06	302	2751.000	2017-01-11	2017-11-08
rf_07	302	2751.000	2017-01-11	2017-11-08
rf_08	492	1705.228	2017-01-11	2018-05-17
rf_09	493	1701.856	2017-01-11	2018-05-17
rf_10	139	43.000	2017-11-03	2018-05-17
rf_12	493	43.000	2017-01-11	2018-05-17
rf_13	68	43.000	2017-01-11	2017-03-18
rf_14	167	43.000	2017-05-04	2018-05-17
rf_15	493	43.000	2017-01-11	2018-05-17
rf_16	493	43.000	2017-01-11	2018-05-17
rf_17	289	43.000	2017-01-18	2018-05-17
rf_18	493	43.000	2017-01-11	2018-05-17
rf_20	493	43.000	2017-01-11	2018-05-17
rf_21	493	43.000	2017-01-11	2018-05-17
rf_22	126	43.000	2017-05-01	2018-05-17
rf_25	492	43.000	2017-01-11	2018-05-17
rf_26	94	43.000	2017-08-21	2017-12-17
rf_27	493	43.000	2017-01-11	2018-05-17
rf_28	493	43.000	2017-01-11	2018-05-17
rf_29	3	43.000	2017-01-25	2017-01-27
rf_30	493	43.000	2017-01-11	2018-05-17
rf_32	493	43.000	2017-01-11	2018-05-17
rf_34	493	43.000	2017-01-11	2018-05-17
rf_35	363	43.000	2017-05-19	2018-05-17
rf_36	49	43.000	2017-12-06	2018-04-26
rf_38	295	43.000	2017-01-11	2018-05-17
rf_39	124	43.000	2017-05-14	2018-03-20
rf_40	493	43.000	2017-01-11	2018-05-17
rf_41	8	43.000	2017-09-19	2018-04-27
rf_42	451	43.000	2017-02-20	2018-05-17
rf_43	493	43.000	2017-01-11	2018-05-17
rf_45	492	43.000	2017-01-11	2018-05-17
rf_46	85	43.000	2018-02-22	2018-05-17
rf_47	493	43.000	2017-01-11	2018-05-17

6 Data quality analysis

Now produce some data quality plots & tables.

6.1 Circuit label checks

The following table shows the number of data files with different circuit labels by household. In theory there should only be one unique list per household and it should be present in every data file.

NB: This table is only legible in the html version of this report because latex does a very bad job of wrapping table cell text. A version is saved in /Volumes/hum-csafe/Research Projects/GREEN Grid/Clean_data/safe/gridSpy/1min/circuitLabelCheck.csv for viewing in e.g. xl.

Bed 2, 2nd Fridge\$2828, Heat Pump\$2826, Hot Water - Controlled\$2825, Incomer - Uncontrolled\$2824, Kitchen, Laundry & Bedroom & Lounge Heat Pumps\$2741, Incomer 1 - All\$2738, Incomer 2 - All\$2737, Kitchen Appliances\$2735, Laundry\$2733, Bedrooms & Lounge\$2602, Heat Pump\$2598, Incomer - All\$2599, Kitchen Appliances\$2601, Laundry & Garage\$2597, Over Cooking Bath tile heat\$1573, Fridge\$1572, Heating\$1576, Hot Water\$1574, Lights\$1577, Mains\$1575
 Downstairs (inc 1 Heat Pump)\$2212, Hot Water - Controlled\$2208, Incomer - Uncontrolled\$2209, Kitchen & Laundry\$2213, Fridge\$2752, Heat Pump & Washing Machine\$2750, Incomer - All\$2748, Kitchen Appliances & Garage\$2753, Lower Bedroom Hallway & Washing Machine\$2683, Hot Water - Controlled\$2679, Incomer 1 - Uncont inc Oven\$2681, Incomer 2 - Uncont inc Heat Pump (x2) & Lounge Power\$4166, Hot Water - Controlled\$4167, Incomer - Uncontrolled\$4168, Kitchen Appliances\$4169, Heat Pump & 2 x Bathroom Heat\$4171, Incomer - All\$4170, Kitchen Power & Heat, Lounge\$4174, Laundry, Garage & 2 Bedrooms Heat Pump & Bedroom 2\$2731, Incomer 1 - Uncont - Inc Hob\$2729, Incomer 2 - Uncont - Inc Oven\$2730, Kitchen Appliances\$4186, Hot Water - Controlled\$4184, Incomer - Uncontrolled\$4181, Laundry\$4185, Lighting\$4186, Heat Pump & Lounge\$2590, Hob\$2589, Hot Water Cpbdr Heater- Cont\$2586, Incomer - Uncontrolled\$2585, Kitchen Appliances\$2586, Heat Pump & Misc\$2107, Hob\$2109, Hot Water - Controlled\$2110, Incomer 1 - Uncontrolled\$2112, Incomer 2 - Uncontrolled\$2113, Heat Pump\$2092, Hot Water - Controlled\$2094, Incomer - Uncontrolled\$2093, Kitchen\$2089, Laundry & 2nd Fridge Freezer\$2148, Hot Water - Controlled\$2150, Incomer 1 - Uncont - inc Hob\$2152, Incomer 2 - Uncont - inc Oven\$2151, Lighting\$2758, Hob & Kitchen Appliances\$2759, Hot Water - Controlled\$2761, Incomer 1 - Uncontrolled \$2763, Incomer 2 - Uncontrolled \$2763, Hob & Kitchen Appliances\$2759, Hot Water - Controlled\$2761, Incomer 1 - Uncontrolled \$2763, Incomer 2 - Uncontrolled \$2763, Heat Pump\$4124, Hot Water - Uncontrolled\$4125, Incomer - Uncontrolled\$4126, Kitchen Appliances\$4121, Laundry, Garage & 2nd Fridge Heat Pump\$4130, Hot Water - Uncontrolled\$4131, Incomer - All\$4132, Kitchen Appliances\$4127, Laundry & Freezer\$4128, Lighting\$4134, Hot Water - Controlled\$4135, Incomer -Uncontrolled\$4136, Kitchen Appliances\$4137, Laundry & Fridge Freezer\$4150, Hot Water - Uncontrolled\$4147, Incomer - All\$4148, Kitchen Appliances\$4145, Lighting\$4149, Washing Machine\$4154, Hot Water - Controlled\$4155, Incomer - Uncontrolled\$4156, Kitchen Appliances\$4151, Laundry \$4152, Lighting\$4160, Hot Water - Controlled\$4158, Incomer - Uncontrolled\$4157, Kitchen Appliances\$4161, Laundry & Garage & 2nd Fridge Heat Pump\$4175, Hot Water - Controlled\$4178, Incomer - Uncontrolled\$4177, Kitchen, Dining & Office\$4179, Laundry, Lounge & 2nd Fridge Heat Pump\$4190, Incomer - All\$4192, Kitchen Appliances\$4187, Laundry\$4188, Lighting\$4189, Oven\$4191, Lighting\$4196, Hot Water - Controlled\$4198, Incomer - All\$4193, Kitchen Appliances\$4195, Laundry\$4194, Lighting\$4195, Heat Pump\$4204, Hot Water - Controlled\$4200, Incomer - All\$4199, Kitchen Appliances\$4201, Laundry\$4202, Lighting\$4202, Heat Pump\$4211, Incomer - All\$4213, Kitchen Appliances\$4210, Laundry, Garage & Guest Bed\$4215, Lighting\$4212, Oven\$4219, Incomer - All\$4221, Kitchen Appliances\$4216, Laundry\$4217, Lighting\$4218, PV & Garage\$4220, Heat Pump\$4223, Hot Water - Uncontrolled\$4224, Incomer - All\$4225, Kitchen Appliances\$4226, Laundry & Garage Freezer\$4232, Heat Pumps (2x) & Power\$4399, Hot Water - Controlled\$4231, Hot Water - Controlled\$4231, Heating\$1633, Hot water\$1636, Kitchen power\$1632, Lights\$1635, Mains\$1634, Range\$1637
 Hob\$3954, Hot Water\$3952, Incomer 1\$3956, Incomer 2\$3955, Laundry & Kitchen Appliances\$3951, Oven\$3953
 Hot Water (2 elements)\$4247, Incomer - Uncontrolled\$4248, Kitchen Appliances\$4244, Lighting & 2 Towel Rail\$4245, Oven\$4248, Hot Water - Controlled (HEMS)\$2081, Incomer - Uncontrolled\$2082, Kitchen, Laundry & Ventilation\$2084, Oven\$2085, PV & Heating\$2102, Incomer - Uncontrolled\$2101, Kitchen\$2104, Laundry, Fridge & Freezer\$2105, Oven & Hob\$2105, Hot Water - Controlled\$2129, Incomer 1 - Uncontrolled\$2128, Incomer 2 - Uncontrolled\$2130, Kitchen Appliances & Ventilation\$2236, Incomer - Uncontrolled\$2237, Kitchen & Laundry\$2234, Lighting\$2232, Oven\$2235, Ventilation\$2248, Incomer - Uncontrolled\$2249, Kitchen\$2246, Laundry, Downstairs & Lounge\$2245, Lighting\$2245, Hot Water - Controlled\$2719, Incomer 1 - Uncont inc Stove\$2718, Incomer 2 - Uncont inc Oven\$2717, Kitchen Appliances\$2717, Hot Water - Controlled\$4144, Incomer - Uncontrolled\$4143, Kitchen Appliances & Heat Pump\$4140, Laundry & Teenagers' Room\$4238, Incomer - All\$4239, Kitchen Appliances\$4234, Laundry & Kitchen\$4235, Lighting\$4236, Oven\$4236, Incomer 1 - All\$2703, Incomer 2 - All\$2704, Kitchen Appliances\$2706, Laundry, Sauna & 2nd Fridge\$2707, Oven\$2705, Spa\$2707, Incomer 1 - Hot Water - Cont\$2626, Incomer 2 - Uncontrolled\$2625, Incomer 3 - Uncontrolled\$2627, Kitchen Appliances & Lounge\$5620, Incomer 2 - inc Bottom Oven\$5621, Kitchen Appliances\$5625, Laundry & Garage\$5624, Lighting\$2726, Incomer 2 - Uncontrolled\$2725, Kitchen Appliances & Laundry\$2722, Microwave\$2721, Oven\$2721,

If this is not the case then this implies that:

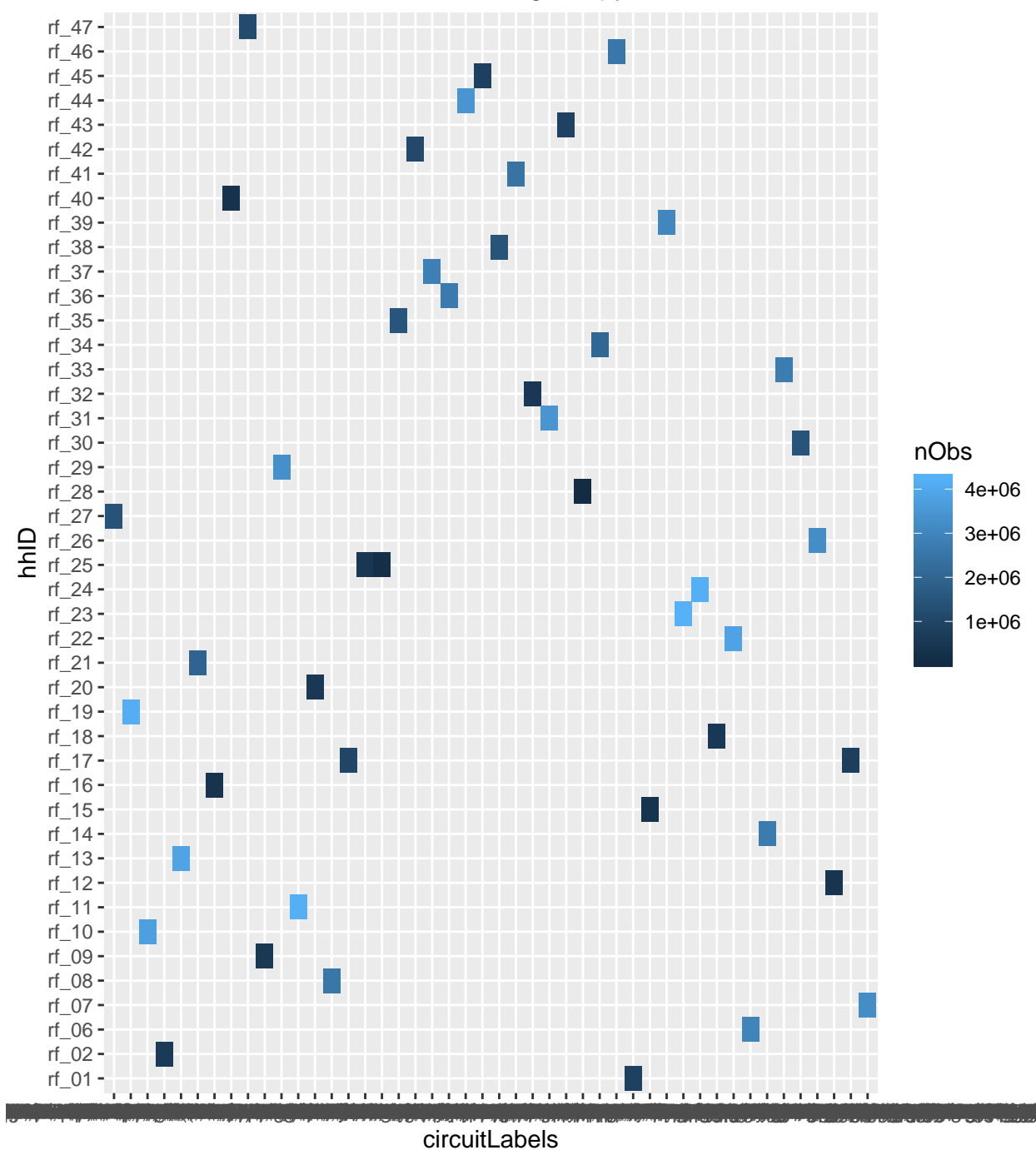
- some of the circuit labels for these households may have been changed during the data collection process;

- some of the circuit labels may have character conversion errors which have changed the labels during the data collection process;
- at least one file from one household has been saved to a folder containing data from a different household (unfortunately the raw data files do *not* contain household IDs in the data or the file names which would enable checking/preventative filtering). This will be visible in the table if two households appear to share *exactly* the same list of circuit labels.

Some or all of these may be true at any given time!

Errors are easy to spot in the following plot where a hhID spans 2 or more circuit labels.

Circuit label counts for all loaded grid spy data



Data source: /Volumes/hum-csafe/Research Projects/GREEN Grid/_RAW DATA/GridSpyData/
Using data received up to 2018-05-18
Only files of size > 3000 bytes loaded

Saving 6.5 x 8 in image

The following table provides more detail to aid error checking. Check for:

- 2+ adjacent rows which have exactly the same circuit labels but different hh_ids. This implies some data from one household has been saved in the wrong folder;
- 2+ adjacent rows which have different circuit labels but identical hh_ids. This could imply the same

thing but is more likely to be errors/changes to the circuit labelling.

If the above plot and this table flag a lot of errors then some re-naming of the circuit labels (column names) may be necessary.

NB: As before, the table is only legible in the html version of this report because latex does a very bad job of wrapping table cell text. A version is saved in /Volumes/hum-csafe/Research Projects/GREEN Grid/Clean_data/safe/gridSpy/1min/circuitLabelMetaDataCheckTable.csv for viewing in e.g. xl.

circuitLabels

Bed 2, 2nd Fridge\$2828, Heat Pump\$2826, Hot Water - Controlled\$2825, Incomer - Uncontrolled\$2824, Kitchen, Laundry & Bedroom & Lounge Heat Pumps\$2741, Incomer 1 - All\$2738, Incomer 2 - All\$2737, Kitchen Appliances\$2735, Laundry\$2734, Bedrooms & Lounge\$2602, Heat Pump\$2598, Incomer - All\$2599, Kitchen Appliances\$2601, Laundry & Garage\$2597, Over Cooking Bath tile heat\$1573, Fridge\$1572, Heating\$1576, Hot Water\$1574, Lights\$1577, Mains\$1575
Downstairs (inc 1 Heat Pump)\$2212, Hot Water - Controlled\$2208, Incomer - Uncontrolled\$2209, Kitchen & Laundry\$2213
Fridge\$2752, Heat Pump & Washing Machine\$2750, Incomer - All\$2748, Kitchen Appliances & Garage\$2753, Lower Bedroom Hallway & Washing Machine\$2683, Hot Water - Controlled\$2679, Incomer 1 - Uncont inc Oven\$2681, Incomer 2 - Uncont inc Heat Pump & 2 x Bathroom Heat\$4171, Incomer - All\$4170, Kitchen Power & Heat, Lounge\$4174, Laundry, Garage & 2 Bedroom Heat Pump & Bedroom 2\$2731, Incomer 1 - Uncont - Inc Hob\$2729, Incomer 2 - Uncont - Inc Oven\$2730, Kitchen Appliances\$4186, Hot Water - Controlled\$4184, Incomer - Uncontrolled\$4181, Laundry\$4185, Lighting\$4187, Heat Pump & Lounge\$2590, Hob\$2589, Hot Water Cpbd Heater- Cont\$2586, Incomer - Uncontrolled\$2585, Kitchen Appliances\$2107, Hob\$2109, Hot Water - Controlled\$2110, Incomer 1 - Uncontrolled\$2112, Incomer 2 - Uncontrolled\$4166, Hot Water - Controlled\$4167, Incomer - Uncontrolled\$4168, Kitchen Appliances\$4169, Heat Pump\$2092, Hot Water - Controlled\$2094, Incomer - Uncontrolled\$2093, Kitchen\$2089, Laundry & 2nd Fridge Freezer\$2148, Hot Water - Controlled\$2150, Incomer 1 - Uncont - inc Hob\$2152, Incomer 2 - Uncont - inc Oven\$2151, Lighting\$2758, Hob & Kitchen Appliances\$2759, Hot Water - Controlled\$2761, Incomer 1 - Uncontrolled \$2763, Incomer 2 - Uncontrolled\$2758, Hob & Kitchen Appliances\$2759, Hot Water - Controlled\$2761, Incomer 1 - Uncontrolled \$2763, Incomer 2 - Uncontrolled\$4124, Hot Water - Uncontrolled\$4125, Incomer - Uncontrolled\$4126, Kitchen Appliances\$4121, Laundry, Garage & 2nd Fridge Heat Pump\$4130, Hot Water - Uncontrolled\$4131, Incomer - All\$4132, Kitchen Appliances\$4127, Laundry & Freezer\$4128, Heat Pump\$4134, Hot Water - Controlled\$4135, Incomer -Uncontrolled\$4136, Kitchen Appliances\$4137, Laundry & Fridge Heat Pump\$4150, Hot Water - Uncontrolled\$4147, Incomer - All\$4148, Kitchen Appliances\$4145, Lighting\$4149, Washing Machine\$4154, Hot Water - Controlled\$4155, Incomer - Uncontrolled\$4156, Kitchen Appliances\$4151, Laundry \$4152, Lighting\$4160, Hot Water - Controlled\$4158, Incomer - Uncontrolled\$4157, Kitchen Appliances\$4161, Laundry & Garage Heat Pump\$4175, Hot Water - Controlled\$4178, Incomer - Uncontrolled\$4177, Kitchen, Dining & Office\$4179, Laundry, Lounge Heat Pump\$4190, Incomer - All\$4192, Kitchen Appliances\$4187, Laundry\$4188, Lighting\$4189, Oven\$4191
Heat Pump\$4196, Hot Water - Controlled\$4198, Incomer - All\$4193, Kitchen Appliances\$4195, Laundry\$4194, Lighting\$4195
Heat Pump\$4204, Hot Water - Controlled\$4200, Incomer - All\$4199, Kitchen Appliances\$4201, Laundry\$4202, Lighting\$4203
Heat Pump\$4211, Incomer - All\$4213, Kitchen Appliances\$4210, Laundry, Garage & Guest Bed\$4215, Lighting\$4212, Oven\$4219, Incomer - All\$4221, Kitchen Appliances\$4216, Laundry\$4217, Lighting\$4218, PV & Garage\$4220
Heat Pump\$4223, Hot Water - Uncontrolled\$4224, Incomer - All\$4225, Kitchen Appliances\$4226, Laundry & Garage Freezer\$4232, Heat Pumps (2x) & Power\$4399, Hot Water - Controlled\$4231, Hot Water - Controlled\$4232, Heating\$1633, Hot water\$1636, Kitchen power\$1632, Lights\$1635, Mains\$1634, Range\$1637
Hob\$3954, Hot Water\$3952, Incomer 1\$3956, Incomer 2\$3955, Laundry & Kitchen Appliances\$3951, Oven\$3953
Hot Water (2 elements)\$4247, Incomer - Uncontrolled\$4248, Kitchen Appliances\$4244, Lighting & 2 Towel Rail\$4245, Oven\$4246, Hot Water - Controlled (HEMS)\$2081, Incomer - Uncontrolled\$2082, Kitchen, Laundry & Ventilation\$2084, Oven\$2085, PV & Ventilation\$2102, Incomer - Uncontrolled\$2101, Kitchen\$2104, Laundry, Fridge & Freezer\$2105, Oven & Hob\$2106, Hot Water - Controlled\$2129, Incomer 1 - Uncontrolled\$2128, Incomer 2 - Uncontrolled\$2130, Kitchen Appliances & Ventilation\$2236, Incomer - Uncontrolled\$2237, Kitchen & Laundry\$2234, Lighting\$2232, Oven\$2235, Ventilation\$2248, Incomer - Uncontrolled\$2249, Kitchen\$2246, Laundry, Downstairs & Lounge\$2245, Lighting\$2246, Hot Water - Controlled\$2719, Incomer 1 - Uncont inc Stove\$2718, Incomer 2 - Uncont inc Oven\$2717, Kitchen Appliances\$2718, Hot Water - Controlled\$4144, Incomer - Uncontrolled\$4143, Kitchen Appliances & Heat Pump\$4140, Laundry & Teenagers\$4141, Hot Water - Controlled\$4238, Incomer - All\$4239, Kitchen Appliances\$4234, Laundry & Kitchen\$4235, Lighting\$4236, Over

circuitLabels

Incomer 1 - All\$2703, Incomer 2 - All\$2704, Kitchen Appliances\$2706, Laundry, Sauna & 2nd Fridge\$2707, Oven\$2705, Spa
Incomer 1 - Hot Water - Cont\$2626, Incomer 2 - Uncontrolled\$2625, Incomer 3 - Uncontrolled\$2627, Kitchen Appliances &
Incomer 1 - Uncontrolled\$2726, Incomer 2 - Uncontrolled\$2725, Kitchen Appliances & Laundry\$2722, Microwave\$2721, Oven
Incomer 1 - inc Top Oven\$5620, Incomer 2 - inc Bottom Oven\$5621, Kitchen Appliances\$5625, Laundry & Garage\$5624, Li

Things to note:

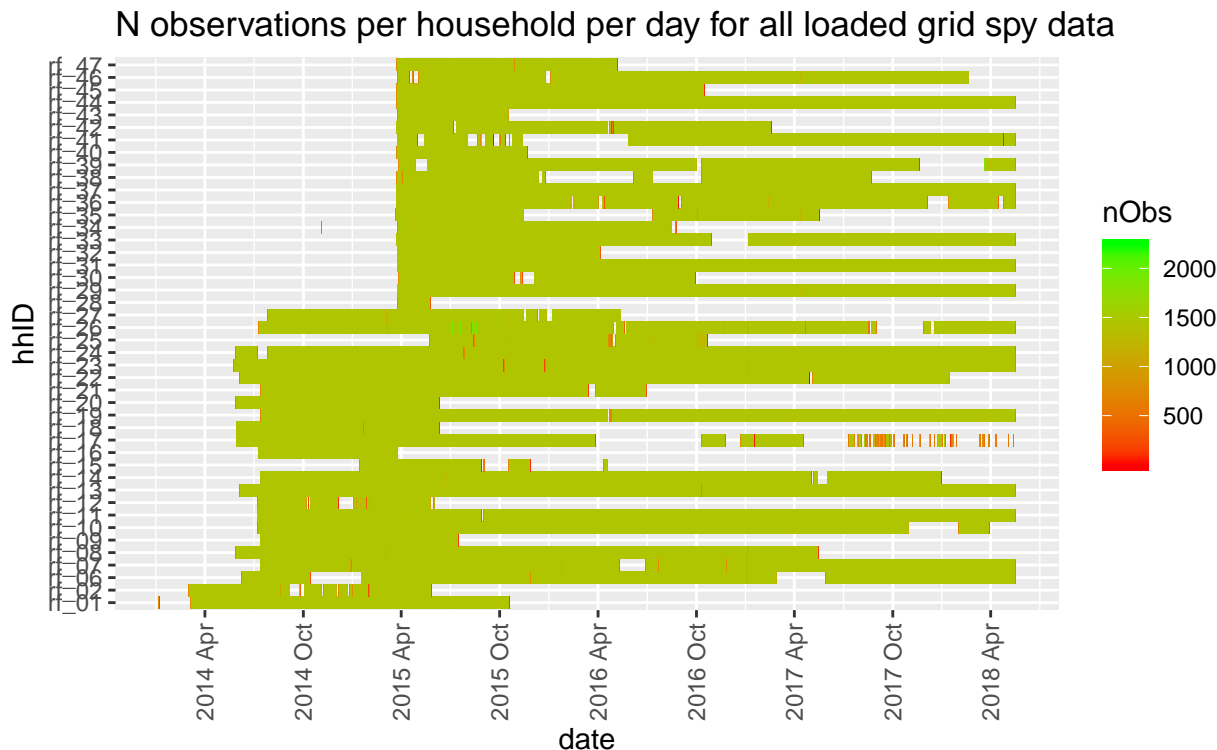
- rf_25 has an additional unexpected “Incomer 1 - Uncontrolled\$2757” circuit in some files but it’s value is always NA

6.2 Observations

The following plots show the number of observations per day per household. In theory we should not see:

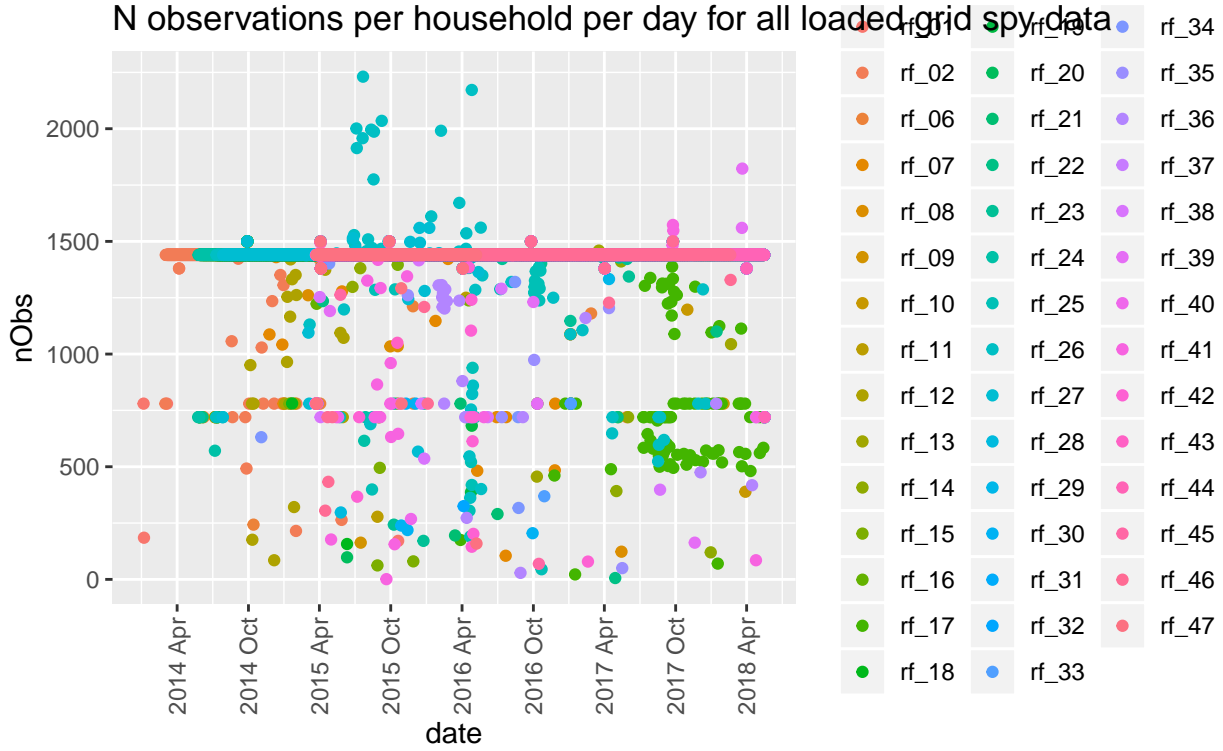
- dates before 2014 or in to the future. These may indicate:
 - date conversion errors;
- more than 1440 observations per day. These may indicate:
 - duplicate time stamps - i.e. they have the same time stamps but different power (W) values or different circuit labels;
 - observations from files that are in the ‘wrong’ rf_XX folder and so are included in the ‘wrong’ household as ‘duplicate’ time stamps.

If present both of the latter may have been implied by the table above and would have evaded the de-duplication filter which simply checks each complete row against all others within it’s consolidated household dataset (a *within household absolute duplicate* check).



Data source: /Volumes/hum-csafe/Research Projects/GREEN Grid/_RAW DATA/GridSpyData/
Using data received up to 2018-05-18
Only files of size > 3000 bytes loaded

Saving 6.5 x 4.5 in image



Saving 6.5 x 4.5 in image

The following table shows the min/max observations per day and min/max dates for each household. As above, we should not see:

- dates before 2014 or in to the future (indicates date conversion errors)
- more than 1440 observations per day (indicates potentially duplicate observations)
- non-integer counts of circuits as it suggests some column errors

We should also not see NA in any row (indicates date conversion errors).

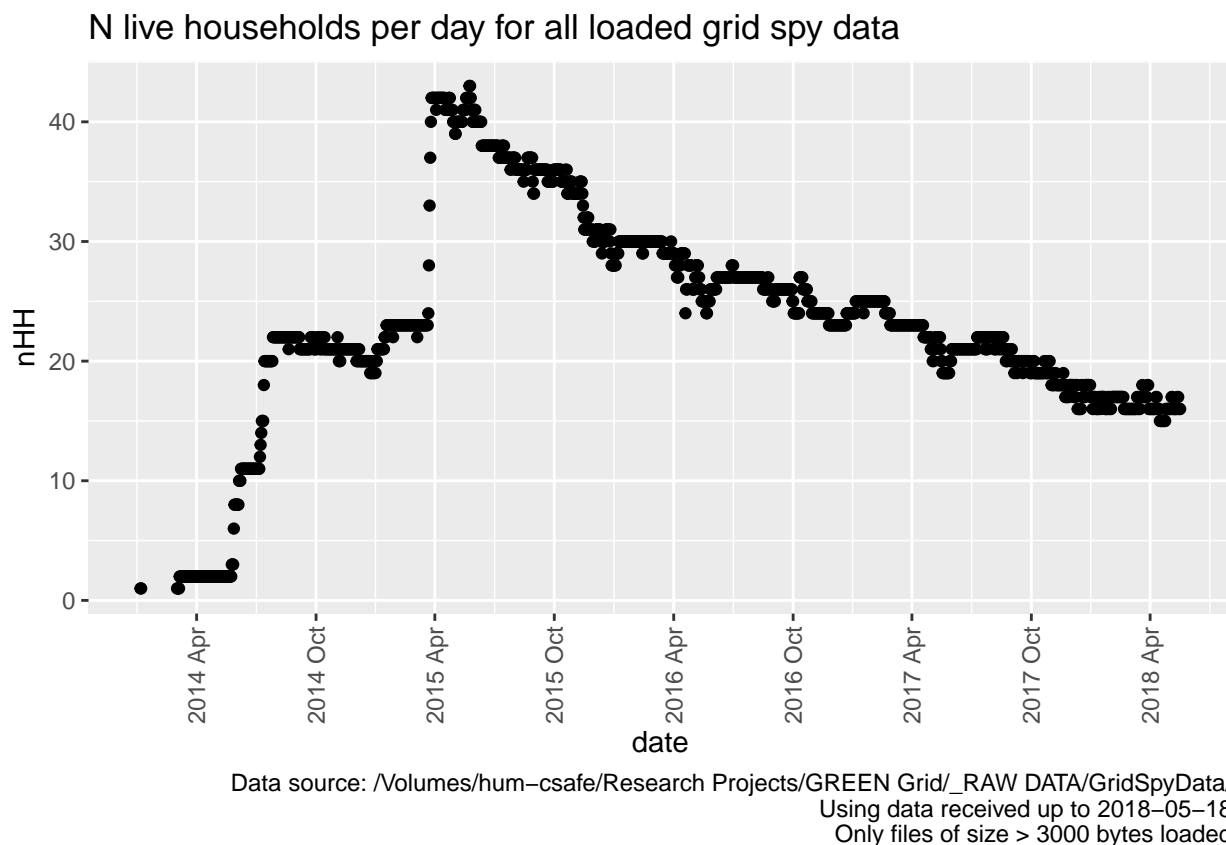
If we do see any of these then we still have data cleaning work to do!

Table 9: Summary observation stats by hhID

hhID	minObs	maxObs	meanNDataColumns	minDate	maxDate
rf_01	171	1500	6	2014-01-05	2015-10-20
rf_02	215	1440	6	2014-03-02	2015-05-28
rf_06	243	1500	6	2014-06-08	2018-05-17
rf_07	105	1500	6	2014-07-13	2018-05-17
rf_08	123	1500	6	2014-05-28	2017-05-15
rf_09	163	1500	6	2014-07-13	2015-07-16
rf_10	389	1500	6	2014-07-08	2018-03-29
rf_11	278	1500	6	2014-07-07	2018-05-17
rf_12	85	1500	6	2014-07-08	2015-06-02
rf_13	456	1500	6	2014-06-05	2018-05-17
rf_14	120	1500	6	2014-07-13	2017-12-30

hhID	minObs	maxObs	meanNDataColumns	minDate	maxDate
rf_15	62	1440	6	2015-01-14	2016-04-18
rf_16	720	1500	6	2014-07-09	2015-03-25
rf_17	22	1500	6	2014-05-29	2018-05-14
rf_18	157	1500	6	2014-05-29	2015-06-11
rf_19	387	1500	9	2014-07-14	2018-05-17
rf_20	98	1500	6	2014-05-28	2015-06-11
rf_21	195	1500	6	2014-07-14	2016-07-01
rf_22	6	1500	6	2014-06-05	2018-01-14
rf_23	171	1500	6	2014-05-25	2018-05-17
rf_24	571	1500	6	2014-05-28	2018-05-17
rf_25	45	1500	6	2015-05-24	2016-10-22
rf_26	362	2231	6	2014-07-10	2018-05-17
rf_27	567	1560	6	2014-07-27	2016-05-13
rf_28	297	1440	6	2015-03-26	2015-05-26
rf_29	720	1500	6	2015-03-25	2018-05-17
rf_30	205	1500	6	2015-03-27	2016-09-29
rf_31	720	1500	6	2015-03-25	2018-05-17
rf_32	325	1500	6	2015-03-25	2016-04-05
rf_33	369	1500	6	2015-03-23	2018-05-17
rf_34	317	1500	6	2014-11-03	2016-08-24
rf_35	50	1500	6	2015-03-22	2017-05-17
rf_36	29	1500	6	2015-03-23	2018-05-17
rf_37	720	1500	6	2015-03-23	2018-05-17
rf_38	398	1500	6	2015-03-24	2017-08-22
rf_39	163	1823	5	2015-03-27	2018-05-17
rf_40	268	1500	6	2015-03-24	2015-11-22
rf_41	1	1573	6	2015-03-25	2018-05-17
rf_42	79	1500	6	2015-03-23	2017-02-18
rf_43	780	1495	6	2015-03-26	2015-10-18
rf_44	720	1500	6	2015-03-24	2018-05-17
rf_45	69	1499	6	2015-03-24	2016-10-15
rf_46	305	1500	13	2015-03-26	2018-02-19
rf_47	159	1500	6	2015-03-24	2016-05-08

Finally we show the total number of households which we think are still sending data.



Saving 6.5 x 4.5 in image

7 Runtime

Analysis completed in 9835.34 seconds (163.92 minutes) using knitr in RStudio with R version 3.4.4 (2018-03-15) running on x86_64-apple-darwin15.6.0.

8 R environment

R packages used:

- base R - for the basics (R Core Team 2016)
- data.table - for fast (big) data handling (Dowle et al. 2015)
- lubridate - date manipulation (Grolemund and Wickham 2011)
- ggplot2 - for slick graphics (Wickham 2009)
- readr - for csv reading/writing (Wickham, Hester, and Francois 2016)
- dplyr - for select and contains (Wickham and Francois 2016)
- progress - for progress bars (Csárdi and FitzJohn 2016)
- knitr - to create this document & neat tables (Xie 2016)
- kableExtra - for extra neat tables (Zhu 2018)
- nzGREENGrid - for local NZ GREEN Grid project utilities

Session info:

R version 3.4.4 (2018-03-15)

```
## Platform: x86_64-apple-darwin15.6.0 (64-bit)
## Running under: macOS High Sierra 10.13.4
##
## Matrix products: default
## BLAS: /Library/Frameworks/R.framework/Versions/3.4/Resources/lib/libRblas.0.dylib
## LAPACK: /Library/Frameworks/R.framework/Versions/3.4/Resources/lib/libRlapack.dylib
##
## locale:
## [1] en_GB.UTF-8/en_GB.UTF-8/en_GB.UTF-8/C/en_GB.UTF-8/en_GB.UTF-8
##
## attached base packages:
## [1] stats      graphics  grDevices  utils      datasets  methods   base
##
## other attached packages:
## [1] progress_1.1.2      dplyr_0.7.4        readr_1.1.1
## [4] lubridate_1.7.4     data.table_1.10.4-3 kableExtra_0.8.0
## [7] knitr_1.20          ggplot2_2.2.1.9000 nzGREENGrid_0.1.0
##
## loaded via a namespace (and not attached):
## [1] Rcpp_0.12.16      highr_0.6          bindr_0.1.1
## [4] pillar_1.2.2      compiler_3.4.4     plyr_1.8.4
## [7] prettyunits_1.0.2 tools_3.4.4        digest_0.6.15
## [10] evaluate_0.10.1   tibble_1.4.2       gtable_0.2.0
## [13] viridisLite_0.3.0 pkgconfig_2.0.1    rlang_0.2.0.9001
## [16] rstudioapi_0.7    yaml_2.1.18        bindrcpp_0.2.2
## [19] withr_2.1.2       stringr_1.3.0      http_1.3.1
## [22] xml2_1.2.0        hms_0.4.2          rprojroot_1.3-2
## [25] grid_3.4.4        glue_1.2.0         R6_2.2.2
## [28] rmarkdown_1.9     magrittr_1.5        backports_1.1.2
## [31] scales_0.5.0.9000 htmltools_0.3.6    assertthat_0.2.0
## [34] rvest_0.3.2       colorspace_1.3-2    labeling_0.3
## [37] stringi_1.1.7     lazyeval_0.2.1     munsell_0.4.3
```

References

- Csárdi, Gábor, and Rich FitzJohn. 2016. *Progress: Terminal Progress Bars*. <https://CRAN.R-project.org/package=progress>.
- Dowle, M, A Srinivasan, T Short, S Lianoglou with contributions from R Saporta, and E Antonyan. 2015. *Data.table: Extension of Data.frame*. <https://CRAN.R-project.org/package=data.table>.
- Grolemund, Garrett, and Hadley Wickham. 2011. “Dates and Times Made Easy with lubridate.” *Journal of Statistical Software* 40 (3): 1–25. <http://www.jstatsoft.org/v40/i03/>.
- R Core Team. 2016. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing. <https://www.R-project.org/>.
- Wickham, Hadley. 2009. *Ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York. <http://ggplot2.org>.
- Wickham, Hadley, and Romain Francois. 2016. *Dplyr: A Grammar of Data Manipulation*. <https://CRAN.R-project.org/package=dplyr>.
- Wickham, Hadley, Jim Hester, and Romain Francois. 2016. *Readr: Read Tabular Data*. <https://CRAN.R-project.org/package=readr>.

R-project.org/package=readr.

Xie, Yihui. 2016. *Knitr: A General-Purpose Package for Dynamic Report Generation in R*. <https://CRAN.R-project.org/package=knitr>.

Zhu, Hao. 2018. *KableExtra: Construct Complex Table with 'Kable' and Pipe Syntax*. <https://CRAN.R-project.org/package=kableExtra>.