

# Data manipulation with Rcell (Version 1.2-5)

Alan Bush

October 7, 2013

## 1 Introduction

Once you have your data loaded into **R**, you can filter it and plot it as shown in “Getting Started with Rcell”. To read that document type in the console

```
> vignette('Rcell')
```

But many times we want to do some manipulation or transformations on the data before plotting it. In this document you’ll see how this can be done using **Rcell**.

## 2 Transforming variables

If you haven’t done so, load the **Rcell** package and the example dataset with

```
> library(Rcell)
> data(ACL394filtered)
```

The easiest way to modify your dataset is to create new variables from existing ones. For example, its desirable to correct the fluorescence measure of a cell by the background fluorescence. To do this for the YFP channel we can use the `f.bg.y` variable, that contains the most common value (mode) for pixels not associated with any cell. If a cell has no fluorophores, we expect it to have a total fluorescence equivalent to `f.bg.y` times the number of pixels of the cell, `a.tot`. So the background corrected fluorescence can be calculated as `f.tot.y - f.bg.y*a.tot`. To create a new variable called `f.total.y` with the corrected value for fluorescence we can use the `transform` function. As all other **Rcell** functions, the first argument is the `cell.data` object to transform.

```
> X<-transform(X, f.total.y=f.tot.y-f.bg.y*a.tot)
```

Once created, you can use the new variable as any other variable of the dataset. You can create several variables in a single call to `transform`, as shown next for the fluorescence density variables.

```
> X<-transform(X, f.density.y=f.tot.y/a.tot, f.density.c=f.tot.c/a.tot)
```

You can keep track of the variables you’ve created with the `summary` function, that will display among other things the “transformed” variables with their definition.

```
> summary(X)
```

| pos | alpha.factor |
|-----|--------------|
| 1   | 1.25         |
| 2   | 1.25         |
| 3   | 1.25         |
| 8   | 2.50         |
| 9   | 2.50         |
| 10  | 2.50         |
| 15  | 5.00         |
| 16  | 5.00         |
| 17  | 5.00         |
| 22  | 10.00        |
| 23  | 10.00        |
| 24  | 10.00        |
| 29  | 20.00        |
| 30  | 20.00        |
| 31  | 20.00        |

Table 1: example data.frame to merge

### 3 Merging variables

Sometimes there is no formula to specify the new variable you want to create. For example, you might want to create a variable that describes the treatment each position received. In the example dataset (`help(ACL394)`) each position received a different dose of alpha-factor pheromone, according to the Table 1.

You can create this table in Excel<sup>1</sup> and save it as a tab delimited text file. If you name it “mytable.txt”, then you can load it into **R** with `read.table`. The best option is to save the file in your working directory, or to change your working directory to where you saved the file (see `?setwd`).

```
> mytable<-read.table("mytable.txt", head=TRUE)
```

If the first row of your text file contains the column names (recommended), you have to set `head` to `TRUE` in `read.table`. Once loaded you can add the new data to your dataset using the `merge` function. This function looks for common variables between `X` and `mytable` and, if it finds them it merges the dataset according to those common variables. Be aware that the names of the columns of `mytable` have to match EXACTLY<sup>2</sup> to the variables of `X`<sup>3</sup>. In this case it will merge by `pos`. You can also specify the variable to merge by with the `by` argument.

```
> X<-merge(X, mytable)
```

```
merging by pos
```

```
merged vars:
```

```
alpha.factor: numeric w/values 1.25, 2.5, 5, 10, 20
```

### 4 Transform By

A common transformation is normalization, i.e. dividing the value of a variable by the “basal” level. For example, we might be interested in the fold increase of YFP fluorescence through time. So we need to divide the measured value at each time by the value at time zero, and we need to do this for every cell. How can we do this? The steps we should follow are the following:

<sup>1</sup>or from R: `mytable<-data.frame(pos=with(X,unique(pos)),alpha.factor=rep(c(1.25,2.5,5,10,20),each=3))`

<sup>2</sup>**R** is case-sensitive so “pos” is different to “Pos”

<sup>3</sup>You can see these variables with `summary(X)`

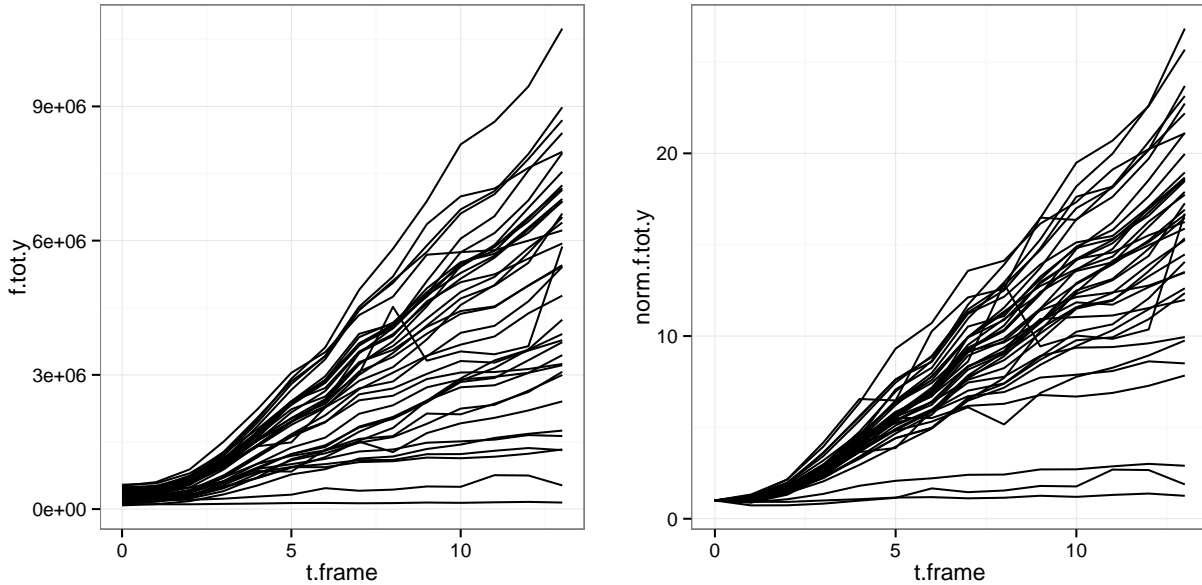


Figure 1: Left: raw single cell time course for YFP fluorescence. Right: Same data normalized to each cells value at time zero.

1. Divide the dataset by cell, creating a table for each cell.
2. Identify the value of fluorescence for time zero.
3. Create a new variable by dividing the fluorescence at each time by the value at time zero.
4. Join the cells datasets back together to retrieve the original dataset with the new variable.

All these steps are done by the function `transformBy`, but it requires information on how each step should be done. For the first step, it needs to know how to partition the dataset. This is specified by passing a quoted list of variable(s) that define the groups. For example, if you want to divide the dataset by position, the second argument of `transformBy` should be `.(pos)`. If you want to divide your dataset by cell use `.(pos,cellID)`. Note that cells in different position can have the same cellID, so the combination of `pos` and `cellID` uniquely identifies a cell. The variable `ucid` (for Unique Cell ID) is another way to uniquely identify a cell. Next we need to specify the name of the new variable to be created (`norm.f.tot.y` for example), and the definition for this variable, `f.tot.y/f.tot.y[t.frame==0]`. With the square brackets we are selecting the value of `f.tot.y` when `t.frame` is zero. Remember to use the logical operator `==` and not the assignment operator `=` within the brackets.

```
> X<-transformBy(X, .(pos,cellID), norm.f.tot.y=f.tot.y/f.tot.y[t.frame==0])
```

You can see the raw and normalized data in Figure 1. Another way to normalize the data, is dividing by the mean of the first three values.

```
> X<-transformBy(X, .(pos,cellID), norm2.f.tot.y=f.tot.y/mean(f.tot.y[t.frame<=2]))
```

## 5 Aggregating your data

To calculate summary statistics you can use the `aggregateBy` function, that returns an aggregated table. That means that the value of each cell of this aggregated table is calculated from more than one cell of the original table. For example, you might be interested in getting the mean YFP fluorescence for each pheromone dose. `aggregateBy` has a similar notation to `transformBy`, where the second argument should be a quoted list of variable names that define the groups by which the aggregation is going to be done. The `select` argument defines which variables are selected for the aggregation.

```
> aggregateBy(X, .(AF.nM), select="f.total.y")
```

```
  AF.nM f.total.y
1  1.25  1071898
2  2.50  1622198
3  5.00  2234246
4 10.00  2393427
5 20.00  2377602
```

You can calculate other statistics using the `FUN` argument, and you can include more than one variable. Here we calculate the median for `f.tot.y`, `f.tot.c` and `a.tot`. Note the use of the wildcard in the `select` argument.

```
> aggregateBy(X, .(AF.nM), select=c("f.tot.*", "a.tot"), FUN=median)
```

```
  AF.nM  f.tot.c f.tot.y a.tot
1  1.25 1047808.5 1212391 415.5
2  2.50 1055751.0 1564543 415.0
3  5.00 1037465.0 2032817 407.0
4 10.00 1001638.5 2224172 398.0
5 20.00  961167.5 2083660 380.0
```

The partition of the dataset can be done by more than one variable, for example by dose and time. Using the function `funstofun` from the `reshape` package, you can calculate more than one statistic at once.

```
> aggregateBy(X, .(t.frame, AF.nM), select="f.density.y", FUN=funstofun(median, sd),
+           subset=t.frame%%6==0)
```

```
  t.frame AF.nM median.f.density.y sd.f.density.y
1         0  1.25          1037.606          162.4781
2         6  1.25          3646.671          892.4902
3        12  1.25          4081.666         1295.7236
4         0  2.50          1043.603          141.9683
5         6  2.50          4663.589         1387.7653
6        12  2.50          6715.928         2200.4825
7         0  5.00          1057.403          163.5038
8         6  5.00          5619.768         1633.4406
9        12  5.00          9600.128         2859.8918
10        0 10.00          1009.415          131.5840
11         6 10.00          5883.817         1218.5057
12        12 10.00         10792.747         2478.7091
13         0 20.00          1022.599          120.6829
14         6 20.00          6440.710         1645.1310
15        12 20.00         11746.938         3383.9809
```

## 6 Evaluating expressions in your dataset

Using the `with` function, you can evaluate an expression in an environment created from your dataset. That means that you can use the names of your variables directly, without any prefix. For example to calculate the mean of `f.tot.y` from position 1

```
> with(X, mean(f.tot.y[pos==1]))
[1] 1372297
```

If you don't use `with` you have to write the full identifier of the variable, and the code becomes longer and harder to understand. For example, the same result can be obtained with

```
> mean(X$data$f.tot.y[X$data$pos==1])
```

## 7 Exporting your data

Although you can do much of your analysis using Rcell functions, you might need to export the data to some other application or use another package within **R**. To retrieve the entire dataset in a `data.frame`, use the double square brackets notation. This returns the registers that pass the QC.filter.

```
> df<-X[[]]
```

This dataset is usually big, and has many variables or registers you are not interested in. You can subset the dataset as you would a `data.frame` (but using double brackets)

```
> df<-X[[pos==1, c("cellID", "f.tot.y", "a.tot")]]
```

You can then save the `data.frame` to a file with `write.table`, or use it in another **R** package.

**Rcell** also provides the function `write.delim`, a small wrapper over `write.table` that creates nice tab delimited text files. For some kinds of data analysis you need your data in a different form than the one **Rcell** uses. You can use the `reshape` function to reshape your data. For instance, a common restructuring is to display time as different columns, and individual cells as different rows. You can obtain this sort of `data.frames` with the following command.

```
> reshape(X, pos+cellID~variable+t.frame, select="f.tot.y",
+         subset=pos<=2&cellID<=5&t.frame%%2==0)
```

|   | pos | cellID | f.tot.y_0 | f.tot.y_2 | f.tot.y_4 | f.tot.y_6 | f.tot.y_8 | f.tot.y_10 | f.tot.y_12 |
|---|-----|--------|-----------|-----------|-----------|-----------|-----------|------------|------------|
| 1 | 1   | 1      | 378752    | 748712    | 1350707   | 2028179   | 2155404   | 2072739    | 2214004    |
| 2 | 1   | 2      | 176429    | 300842    | 448582    | 535334    | 549019    | 562208     | 512430     |
| 3 | 1   | 3      | 384393    | 665472    | 1234888   | 1913377   | 2036718   | 2217148    | 2071306    |
| 4 | 1   | 4      | 245876    | 510412    | 887509    | 1493615   | 1692185   | 1987466    | 2137951    |
| 5 | 2   | 2      | 387551    | 620656    | 1049458   | 1327046   | 1317000   | 1409672    | 1742833    |
| 6 | 2   | 3      | 428014    | 655421    | 1239405   | 1616264   | 1942105   | 2268256    | 2583064    |
| 7 | 2   | 4      | 452047    | 718126    | 1381880   | 1808801   | 2260311   | 2624726    | 3127893    |
| 8 | 2   | 5      | 330852    | 228660    | 285463    | 294746    | 357041    | 354951     | 308372     |

see `help(reshape.cell.data)` for more details.