

eulerr under the hood

Johan larsson

November 15, 2017

1 Introduction

eulerr relies on an extensive machinery to turn user input into a pretty Euler diagram. Little of this requires any tinkering from the user. To make that happen, however, **eulerr** needs to make several well-formed decisions about the design of the diagram on behalf of the user, which is not a trivial task.

This document outlines the implementation of **eulerr** from input to output. It is designed to an evolving documentation on the innards of the program.

2 Input

The workhorse of **eulerr** is `euler()`. To start with, we need input in the form of

- a named numeric vector, such as `c(A = 10, B = 5, "A&B" = 3)`, where ampersands define disjoint set combinations or unions, depending on the argument input,
- a `data.frame` or matrix of logicals or binary indices where each row denotes the set relationships of

- either a single observation

```
matrix(sample(c(TRUE, FALSE), 12, replace = TRUE),
        ncol = 3,
        dimnames = list(NULL, c("A", "B", "C")))
```

```
##           A      B      C
## [1,]  TRUE  TRUE  TRUE
## [2,]  TRUE FALSE  TRUE
## [3,]  TRUE FALSE  TRUE
## [4,] FALSE  TRUE  TRUE
```

- or of a unique set combination if a numeric vector is supplied to the argument `weights`,

```
matrix(c(TRUE, FALSE, FALSE,
          TRUE, TRUE, FALSE,
          FALSE, FALSE, TRUE),
        ncol = 3,
        dimnames = list(NULL, c("A", "B", "C")))
```

```
##           A      B      C
## [1,]  TRUE  TRUE FALSE
## [2,] FALSE  TRUE FALSE
## [3,] FALSE FALSE  TRUE
```

3 Pre-processing

- a table (max 3 dimensions),

```
as.table(apply(Titanic, 2:4, sum))

## , , Survived = No
##
##      Age
## Sex      Child Adult
## Male      35  1329
## Female    17   109
##
## , , Survived = Yes
##
##      Age
## Sex      Child Adult
## Male      29   338
## Female    28   316
```

- or a list of sample spaces, such as

```
list(A = c("x", "xy", "xyz"),
     B = c("xy"),
     C = c("x", "xyz"))

## $A
## [1] "x"  "xy" "xyz"
##
## $B
## [1] "xy"
##
## $C
## [1] "x"  "xyz"
```

If the `data.frame` or `matrix` form is used, the user additionally has the option to split the data set by a factor and compute separate euler diagrams for each split. This is accomplished by supplying a factor variable to the `by` arguments (see the documentation in `?base::by`).

3 Pre-processing

eulerr organizes the input of the user into a matrix of binary indexes, which in **R** is represented as a matrix of logicals. For a three set configuration, this looks like this,

```
library(eulerr)
eulerr:::bit_indexr(3)

##      [,1] [,2] [,3]
## [1,] TRUE FALSE FALSE
## [2,] FALSE TRUE FALSE
```

4 Initial configuration

```
## [3,] FALSE FALSE TRUE
## [4,] TRUE TRUE FALSE
## [5,] TRUE FALSE TRUE
## [6,] FALSE TRUE TRUE
## [7,] TRUE TRUE TRUE
```

and is accompanied by a vector of the *disjoint* areas of the set combinations.

To provide a starting configuration, we work exclusively with circles and, given these areas, we figure out the required pairwise distance between the sets to achieve a circle–circle overlap that matches the set intersection between the sets. We do this numerically, using the formula for a circle–circle overlap,

$$A = r_1^2 \arccos\left(\frac{d^2 + r_1^2 - r_2^2}{2dr_1}\right) + r_2^2 \arccos\left(\frac{d^2 + r_2^2 - r_1^2}{2dr_2}\right) - \frac{1}{2}\sqrt{(-d + r_1 + r_2)(d + r_1 - r_2)(d - r_1 + r_2)(d + r_1 + r_2)}, \quad (1)$$

where r_1 and r_2 are the radii of the first and second circles respectively and d the distance between the circles.

r_1 and r_2 are known but because d is not, we approximate it using one-dimensional numerical optimization. Our loss function is the squared difference between A and the desired overlap, which we then optimize using **R**'s `optimize()`, which is a “combination of golden section search and successive parabolic interpolation”.

```
r1 <- 0.7 #radius of set 1
r2 <- 0.9 #radius of set 2
overlap <- 1 #area of overlap

stats::optimize(eulerr:::discdisc, #computes the squared loss
  interval = c(abs(r1 - r2), sum(r1, r2)),
  r1 = r1,
  r2 = r2,
  overlap = overlap)

## $minimum
## [1] 0.634
##
## $objective
## [1] 2.985e-10

# minimum is our required distance
```

Now that we have the distances, we can proceed to the next step: computing an initial configuration.

4 Initial configuration

Our initial layout can be setup in a number of ways; **eulerr** uses one of the methods from Fredrickson's [venn.js](#), which features a

constrained version of multi-dimensional scaling (MDS) based on that of Wilkinson's R package `venneuler` Wilkinson [1]. `venneuler` tries to place disjoint and subset exactly neck-in-neck and at the exact midpoint of the set respectively. However, since we are indifferent about where in the space outside (or respectively inside) the sets are placed, that behavior becomes problematic since it might interfere with locations of other sets that need to occupy some of that space.

The MDS algorithm from `venn.js` circumvents this by assigning a loss and gradient of 0 when, for instance, the set relationships *and* the candidate ellipses are disjoint. Then, to optimize the pairwise relationships between sets, `eulerr` uses the following loss, gradient, and Hessian functions:

$$\mathcal{L}(h, k) = \sum_{0 \leq i < j \leq N} \begin{cases} 0 & F_i \cap F_j = \emptyset \text{ and } O_{ij} = \emptyset \\ 0 & (F_i \subseteq F_j \text{ or } F_i \supseteq F_j) \text{ and } O_{ij} = \emptyset \\ \left((h_i - h_j)^2 + (k_i - k_j)^2 - d_{ij}^2 \right)^2 & \text{otherwise} \end{cases} \quad (2)$$

The analytical gradient (3) is retrieved as usual by taking the derivative of the loss function:

$$\vec{\nabla} f(h_i) = \sum_{j=1}^N \begin{cases} \vec{0} & F_i \cap F_j = \emptyset \text{ and } O_{ij} = \emptyset \\ \vec{0} & (F_i \subseteq F_j \text{ or } F_i \supseteq F_j) \text{ and } O_{ij} = \emptyset \\ 4(h_i - h_j) \left((h_i - h_j)^2 + (k_i - k_j)^2 - d_{ij}^2 \right) & \text{otherwise,} \end{cases} \quad (3)$$

where $\vec{\nabla} f(k_i)$ is found as in (3) with h_i swapped for k_i (and vice versa), F are the various sets in the input, O are the pairwise overlaps between circles. Because it speeds up convergence, we also compute the Hessian matrix (4). (In our implementation, we only actually use the lower triangle.)

$$H = \sum_{0 \leq i < j \leq N} \begin{bmatrix} 4\left((h_i - h_j)^2 + (k_i - k_j)^2 - d_{ij}^2\right) + 8(h_i - h_j)^2 & \cdots & 8(h_i - h_j)(k_i - k_j) \\ \vdots & \ddots & \vdots \\ 8(k_i - k_j)(h_i - h_j) & \cdots & 4\left((h_i - h_j)^2 + (k_i - k_j)^2 - d_{ij}^2\right) + 8(k_i - k_j)^2 \end{bmatrix}. \quad (4)$$

Note that the constraints given in (2) and (3) still apply to each element of (4) and have been omitted for practical reasons only.

Fredrickson uses the *Polak–Ribière Conjugate Gradient Method* to optimize the initial layout. In our experience, this method occasionally encounters local minima, which is why we have opted to use `nlmminb()` from the R core package `stats`, which is a translation from FORTRAN code developed by Gray [2] and ported to R by Douglas Bates and Deepayan Sarkar, and, although it is a piece of complicated code, performs well for the difficult problem of aligning Euler diagrams.

This initial configuration will work perfectly for any 1–2 set combinations and as well as possible with 3 sets if we use circles but for all other combinations there is usually a need to fine tune the configuration.

5 Final configuration

In order to finalize the configuration we need to be able to compute the areas of the overlaps of the sets, which as it turns out, is *not* trivial. In fact, most of methods rely on approximations of the areas by, for instance, quad-tree binning (**venneuler**) or polygon intersections (**VennMaster** [3]). These methods yield reasonable estimates but, given that the computation may have to run for a vast number of iterations, are usually prohibitive in terms of performance.

venn.js and **eulerAPE** both, however, use exact algorithms. Based on the fact that any intersection of ellipses can be represented as a convex polygon with elliptical segments on the fringes, it is possible to arrive at exact area calculations.

5.1 Intersections

Finding the areas of the overlaps exactly requires that we first know the points at which the different ellipses intersect. **eulerr**'s approach to this is based on a method outlined by [4]. **eulerr** owes significant debt to the R package **RConics** [5], which has been tremendously helpful in developing and, especially, debugging the algorithm. Some parts of the code are in fact straight-up translations to C++ from the code in **RConics**.

The method is based in *projective geometry* (rather than euclidean). To find the intersection points, the algorithm first

- converts the two ellipses from canonical form to matrix notation. The canonical form of a rotated ellipse is given by

$$\frac{((x - h) \cos(\phi) + (y - k) \sin(\phi))^2}{a^2} + \frac{((x - h) \sin(\phi) - (y - k) \cos(\phi))^2}{b^2} = 1,$$

where ϕ is the counter-clockwise angle from the positive x axis to the semi-major axis a . b is the semi-minor axis whilst (h, k) is the center of the ellipse. This is then converted to the matrix form

$$E = \begin{bmatrix} A & B/2 & D/2 \\ B/2 & C & E/2 \\ D/2 & E/2 & F \end{bmatrix},$$

which may be used to represent any conic. We then

- split one of the ellipses (conics) into a pencil of two lines, and subsequently
- intersect the remaining conic with these two lines, which will yield between 0 and 4 intersection points.

5.2 Areas

The next step is to calculate the area of overlap between all the possible combinations of ellipses. A solution for this, albeit only for circles, was first published by Fredrickson in a [blog post](#). It relies on finding all the intersection points between the currently examined sets that

are also within these sets. It is then trivial to find the area of the convex polygon that these vertices make up. Finding the rest of the area, which is made up of the ellipse segments between subsequent points, requires a bit of trigonometry.

Here, we have used an algorithm from [?], which computes circle integral between the points on the ellipse minus the area of the triangle made up of the center of the ellipse:

$$A(\theta_0, \theta_1) = F(\theta_1) - F(\theta_0) - \frac{1}{2}|x_1y_0 - x_0y_1|,$$

$$\text{where } F(\theta) = \frac{a}{b} \left[\theta - \arctan \left(\frac{(b-a) \sin 2\theta}{b+a+(b-a) \cos 2\theta} \right) \right]$$

As our loss function, we use the sum of squared differences between the disjoint set intersections and the areas we have computed and again use the `nlm()` optimizer to layout the set.

In *rare* instances, the algorithm breaks down because of numerical precision errors that might, for instance, lead to some intersection points being left out. In these cases, we rely instead on a approximation by sampling points inside the ellipses and computing the area of the required overlap as the proportion of points inside that overlap to those in the ellipse multiplied with the ellipse's area.

This optimization step is the bottleneck of the overall computations in terms of performance, being that we're optimizing over five parameters for every ellipse (or 3 in the case of circles)—nevertheless, we're profiting from the implementation in the C++ programming language through **Rcpp** [6] and its plugin for the linear algebra library **Armadillo** [7]. When the number of sets is low, **eulerr** outperforms **venneuler**, whereas the relationship reverses as the number of sets increase (around seven), when the sheer number of overlaps that **eulerr** has to examine bogs it down.

In a future version, it is possible that a alternative, approximative method will be introduced to deal with relationships with large numbers of sets.

6 Layout

Since the optimization steps are unconstrained, we run the risk of ending up with dispersed layouts. To fix this, we use a SKYLINE-BL rectangle packing algorithm [8] to pack the disjoint clusters of ellipses (in case there are any) into a heuristically chosen bin.

At the time of writing this algorithm is crudely implemented – for instance, it does not attempt to rotate the rectangles (boundaries for the ellipses) or attempt to use. Since we're dealing with a rather simple version of the rectangle packing problem, however, it seems to do the trick.

7 Output

Before we get to plotting the solution, it is useful to know how well the fit from **eulerr** matches the input. Sometimes euler diagrams are

7 Output

just not feasible, particular for combinations with many sets, in which case we should stop here and look for another design to visualize the set relationships.

It is not, however, obvious what it means for a euler diagram to “fit well”. **venneuler** uses a metric called *stress*, which is defined as

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

where \hat{y}_i is an ordinary least squares estimate from the regression of the fitted areas on the original areas that is being explored during optimization.

Meanwhile, **eulerAPE** [9] uses *diagError*:

$$\max_{i=1,2,\dots,n} \left| \frac{y_i}{\sum y_i} - \frac{\hat{y}_i}{\sum \hat{y}_i} \right|$$

Both metrics are given the user after the diagram has been fit, together with a table of residuals.

```
combo <- c("A" = 1, "B" = 1, "C" = 1,
           "A&B" = 0.5, "A&C" = 0.5, "C&B" = 0.5)

fit1 <- euler(combo)
fit1
```

##	original	fitted	residuals	regionError
## A	1.0	1.038	-0.038	0.021
## B	1.0	1.038	-0.038	0.021
## C	1.0	1.038	-0.038	0.021
## A&B	0.5	0.302	0.198	0.040
## A&C	0.5	0.302	0.198	0.040
## B&C	0.5	0.302	0.198	0.040
## A&B&C	0.0	0.247	-0.247	0.058
##				
## diagError:	0.058			
## stress:	0.049			

It is clear that this is not a good fit, which we can find out just by looking at the plot (Figure 1). This is a good example of when ellipses come in handy.

```
fit2 <- euler(combo, shape = "ellipse")
fit2
```

##	original	fitted	residuals	regionError
## A	1.0	1.0	0	0
## B	1.0	1.0	0	0
## C	1.0	1.0	0	0
## A&B	0.5	0.5	0	0
## A&C	0.5	0.5	0	0
## B&C	0.5	0.5	0	0

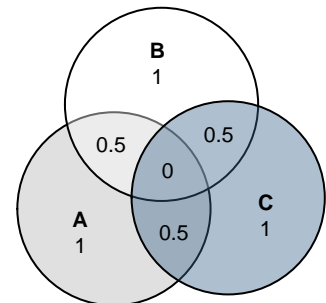


Figure 1. A plot with circles.

```
## A&B&C      0.0      0.0      0      0
##
## diagError: 0
## stress:    0
```

Much better (Figure 2).

8 Plotting

Let's face it: euler diagrams are naught without visualization. Here, **eulerr** interfaces the elegant Lattice graphics system [10] to grant the user extensive control over the output, and allow for faceted plots in case such a design was used in fitting the Euler configuration.

8.1 Labelling

Most users will want to label their Euler diagrams. One option is to simply add a legend

```
plot(euler(c(A = 2, B = 3, "A&B" = 1)), auto.key = TRUE)
```

but many will want to label their diagrams directly, perhaps also adding counts.

```
plot(euler(c(A = 2, B = 3, "A&B" = 1)), counts = TRUE)
```

In this case, laying out the diagram becomes considerably more involved. Finding a reasonable spot for the text inside the diagram only lends itself to an easy solution if the shape of the intersection has a center-of-gravity inside ellipse, in which case an average of some of the points might suffice. This is often not the case, however, and we need a better solution. Specifically, what we need is a method to find the point inside the circle overlap for the counts and circle complement to the intersection for our labels.

So far, we have not been able to derive an analytical solution for finding a good point, or for that matter a reliable way of finding *any* point that is in the required intersection or complement. As is often the case, the next-best thing turns out to be a numerical one. First, we locate a point that is inside the required region by spreading points across one of the discs involved in the set combination. To spread points uniformly, we use *Vogel's method* [11, 12]

$$\left(p_k = (\rho_k, \theta_k) = \left(r \sqrt{\frac{k}{n}}, \pi(3 - \sqrt{5})(k - 1) \right) \right)_{k=1}^n,$$

which is actually based on the golden angle.

```
n <- 500
seqn <- seq(0, n, 1)
theta <- seqn*pi*(3 - sqrt(5))
rad <- sqrt(seqn/n)
x <- rad*cos(theta)
y <- rad*sin(theta)
```

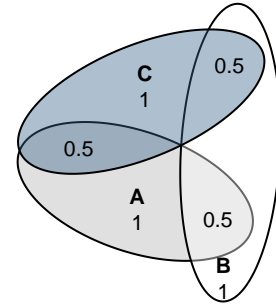


Figure 2. A plot with ellipses.

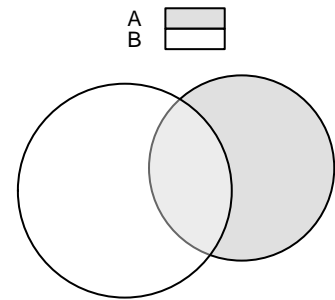


Figure 3. A simple plot with a legend

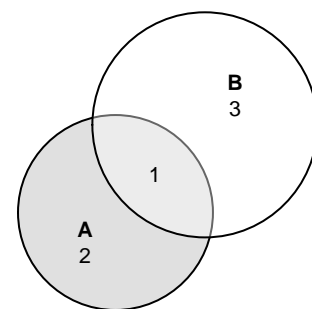


Figure 4. A plot with counts.

After this, we scale, translate, and rotate the points so that they fit the desired ellipse.

After we've spread our points throughout the ellipse and found one that matches our desired combination of ellipses/sets, we then proceed to optimize its position numerically. For this, we use version of the *Nelder–Mead Method* [13] which we've translated from Matlab code by Kelley [14] and customized for **eulerr** (in particular to make sure that the simplex does not escape the intersection boundaries since we for this problem *want* the local minimum).

8.2 Coloring

Per default, the ellipses are filled with colors. The default option is to use an adaptive scheme in which colors are chosen to provide a balance between distinctiveness, beauty, and consideration for the color deficient. The color palette has been generated from **qualpalr** (developed by the author), which automatically generates qualitative color palettes based on a model of color perception.

References

- [1] L. Wilkinson. Exact and approximate area-proportional circular Venn and Euler diagrams. *IEEE Transactions on Visualization and Computer Graphics*, 18(2):321–331, February 2012. ISSN 1077-2626. doi: 10.1109/TVCG.2011.56.
- [2] David M. Gray. Usage summary for selected optimization routines. Technical Report 153, AT&T Bell Laboratories, Murray Hill, NJ, October 1990. URL <https://ms.mcmaster.ca/~bolker/misc/port.pdf>.
- [3] Hans A. Kestler, André Müller, Johann M. Kraus, Malte Buchholz, Thomas M. Gress, Hongfang Liu, David W. Kane, Barry R. Zeeberg, and John N. Weinstein. VennMaster: area-proportional Euler diagrams for functional GO analysis of microarrays. *BMC Bioinformatics*, 9:67, January 2008. ISSN 1471-2105. doi: 10.1186/1471-2105-9-67. URL <https://doi.org/10.1186/1471-2105-9-67>.
- [4] Jürgen Richter-Gebert. *Perspectives on Projective Geometry: A Guided Tour Through Real and Complex Geometry*. Springer, Berlin, Germany, 1 edition, February 2011. ISBN 978-3-642-17286-1.
- [5] Emanuel Huber. RConics: computations on conics, December 2014. URL <https://CRAN.R-project.org/package=RConics>.
- [6] Dirk Eddelbuettel and Romain François. Rcpp: Seamless R and C++ integration. *Journal of Statistical Software*, 40(8):1–18, 2011. URL <http://www.jstatsoft.org/v40/i08/>.

References

- [7] Dirk Eddelbuettel and Conrad Sanderson. RcppArmadillo: accelerating R with high-performance C++ linear algebra. *Computational Statistics and Data Analysis*, 71:1054–1063, March 2014. URL <http://dx.doi.org/10.1016/j.csda.2013.02.005>.
- [8] Jukka Jylänki. A thousand ways to pack the bin – a practical approach to two-dimensional rectangle bin packing, February 2010. URL <http://clb.demon.fi/files/RectangleBinPack.pdf>.
- [9] Luana Micallef and Peter Rodgers. eulerAPE: drawing area-proportional 3-Venn diagrams using ellipses. *PLOS ONE*, 9(7): e101717, 2014-jul-17. ISSN 1932-6203. doi: 10.1371/journal.pone.0101717.
- [10] Deepayan Sarkar. *Lattice: multivariate data visualization with R. Use R!* Springer, New York, USA, 2008. ISBN 978-0-387-75968-5. URL <http://www.springer.com/us/book/9780387759685>.
- [11] Mary K. Arthur. Point picking and distributing on the disc the sphere. Final ARL-TR-7333, US Army Research Laboratory, Weapons and Materials Research Directorate, Aberdeen, USA, July 2015. URL www.dtic.mil/get-tr-doc/pdf?AD=ADA626479.
- [12] H. Vogel. A better way to construct the sunflower head. *Mathematical Biosciences*, 44(3-4):179–189, 1979. doi: 10.1016/0025-5564(79)90080-4.
- [13] J. A. Nelder and R. Mead. A simplex method for function minimization. *The Computer Journal*, 7(4):308–313, January 1965. ISSN 0010-4620. doi: 10.1093/comjnl/7.4.308. URL <https://academic.oup.com/comjnl/article/7/4/308/354237/A-Simplex-Method-for-Function-Minimization>.
- [14] C. T. Kelley. *Iterative methods for optimization*. Number 18 in Frontiers in applied mathematics. Society for Industrial and Applied Mathematics, Philadelphia, USA, 1 edition edition, 1999. ISBN 0-89871-433-8.