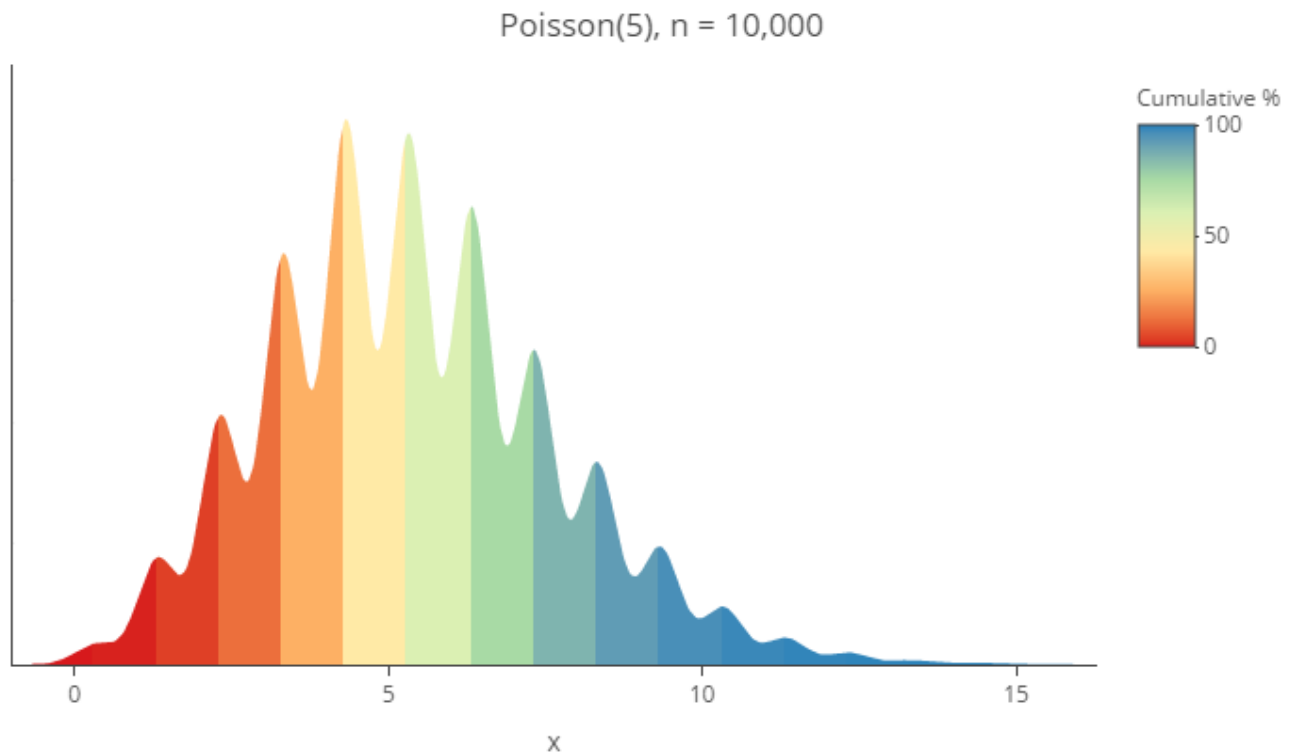# Using Heatmap Coloring on a Density Plot Using R to Visualize Distributions

🌐 **displayr.com**/using-heatmap-coloring-density-plot-using-r-visualize-distributions/



Lots of different visualizations have been proposed for understanding distributions. In this post, I am going show how to create my current favorite, which is a density plot using *heatmap shading*. I find them particularly useful for showing duration data.

## Example: time to purchase

The table below shows the number of days it took 200 people to purchase a product, after their first trial of the product. The average time taken is 131 days and the median is 54 days. The rest of the post looks at a number of different ways of visualizing this data, finishing with the *heated density plot*.

## Density plot

One of the classic ways of plotting this type of data is as a *density plot.* The standard R version is shown below. I have set the default from argument to better display this data, as otherwise density plots tend to show negative values even when all the data contains no negative values.

```
1    (days, from =
   plot(density0),
2
3            main  "Density
             =     plot"              ,

             xlab  "Number of days since trial
             =     started"                    )
```

This plot clearly shows that purchases occur at relatively close to 0 days since the trial started. But, unless you use a ruler, there is no way to work out precisely where the peak occurs. Is it at 20 days or 50 days? The values shown on the *y*-axis have no obvious meaning (unless you read the technical documentation, and even then they are not numbers that can be explained in a useful way to non-technical people).

We can make this easier to read by only plotting the data for the first 180 days ( to = 180), and changing the bandwidth used in estimating the density (adjust = 0.2) to make the plot less smooth. We can now see that the peak is around 30. While the plot does a good job at describing the shape of the data, it does not allow us to draw conclusions regarding the cumulative proportion of people to purchase by a specific time. For example, what proportion of people buy in the first 100 days? We need a different plot.

```
1    (days, from = 0, to = 180, adjust =
   plot(density0.2),
2
3            main  "Density plot - Up to 180 days (86% of
             =     data)"                                  ,

             xlab  "Number of days since trial
             =     started"                    )
```

## Survival curve

*Survival curves* have been invented for this type of data. An example is shown below. In this case, the survival curve shows the proportion of people who have yet to purchase at each point in time. The relatively sharp drop in the *survival function* at 30 is the same pattern that was revealed by the peak in the density plot. We can see also see that at about 100 days, around 26% or so of people have "survived" (i.e., not purchased). The plot also reveals that around 14% of customers have not purchased during the time interval shown (as this is where the line crosses the right-side of the plot). To my mind the survival plot is more informative than the density plot. But, it is also harder work. It is hard to see most non-technical audiences finding all the key patterns in this plot without guidance.

```
1   library(survival)

2   surv.days
    =          Surv(days)
3

    surv.fit
4   =          survfit(surv.days~1)

5       (surv.fit, main
6   plot=
    "Kaplan-Meier estimate with 95% confidence bounds (86% of
7   data)"                                                          ,

8       xlab   "Days since trial
        =      started"                      ,

        xlim   (0,
        =      c180),

        ylab
        =      "Survival function")

        (20, 10, lwd =
    grid2)
```

## Heated density plot

I call the visualization below a *heated density plot.* No doubt somebody invented this before we did, so please tell me if there is a more appropriate name. It is identical to the density plot from earlier in this post, except that:

- The heatmap coloring shows the cumulative proportion to purchase, ranging from red (0%), to yellow (50%, the median), to blue (100%). Thus, we can now see that the median is at about 55, which could not be ascertained from the earlier density plots.

- If you hover your mouse pointer (if you have one) over the plot, it will show you the cumulative percentage of people to purchase. We can readily see that the peak occurs near 30 days. Similarly, we can see that 93% of people purchased within 500 days.

```
1                    (days, from =
    HeatedDensityPlot0,
2
                     title   "Heated density
3                    =       plot"                   ,

4                    x.title   "Number of days since trial
                     =         started"                       ,

                     legend.title   "% of
                     =                 buyers"        )
```

Below I show the plot limited to the first 180 days with the bandwidth adjusted as in the earlier density plot. The use of the legend on the right allows the viewer to readily see that only a little over 85% of people have purchased.

```
1   HeatedDensityPlot0.2,         (days, from = 0, to = 180, adjust =
2
3                     title    "Heated density plot - Up to 180 days (86% of
                      =        data)"                                          ,
4                     x.title   "Number of days since trial
                      =         started"                                 ,
                      legend.title   "% of
                      =              buyers"          )
```

## Bonus: they are awesome for discrete data

The visualization below shows randomly-generated data where I have generated whole numbers in the range of 0 through 12 (i.e., one 0; six 1s, seven 2s, etc). In all the previous heated density plots the color transition was relatively smooth. In the visualization below, the discrete nature of the data has been emphasized by the clear vertical lines between each color. This is not a feature that can be seen on either a traditional density plot or histogram with small sample sizes.

```
1   set.seed(12230)
2   x
    =   rpois(100, 5)
3
4   HeatedDensityPlot0,           (x, from =
5
                      title    "Poisson(5), n =
                      =        100"                      ,
                      x.title
                      =          "x")
```

## Conclusion

The *heated density plot* is, to my mind, a straightforward improvement on traditional density plots. It shows additional information about the cumulative values of the distribution and reveals when the data is discrete. However, the function we have written is still a bit rough around the edges. For example, sometimes a faint line appears at the top of the plot, and a whole host of warnings appear in R when the plot is created. If you can improve on this, please tell us how!

## The code

The R function used to create the plot is shown below. If you run this code you will get a plot of a normal distribution. If you want to reproduce my exact examples, please click here to sign into Displayr and access the document that contains all my code and charts. This document contains further additional examples of the use of this plot. To see the code, just click on one of the charts in the document, and the code should be shown on the right in **Properties > R CODE**.

```
1   HeatedDensityPlot
    <-                    function(x,
2
```

```
 3                              title =          "",
 4
 5                              x.title =          "x",
 6
 7                              colors =          c('#d7191c','#fdae61','#ffffbf','#abdda4',
                                '#2b83ba'),
 8
 9                              show.legend =              TRUE,
10                              legend.title  "Cumulative
11                              =            %"                ,
12                              n = 512,
13
14                              n.breaks = 5,
15                              ...)
16 {
17
18     n.obs <-      length(x)
19
20     if(any<-     is.na{        (nas         (x)))
21
22         warning(sum(nas),removed." " observations with missing values have been                              )
23         x <- x[!nas]
24
25         n.obs <-      length(x)
26
27     }
28     if4)   (n.obs <
29         stop(n.obs,
" is too few observations for a valid density plot. See Silverman (1986)
Density Estimation for Statistics and Data Analysis."
30 )
31
32     dens           (x,
33     <-      density...)
34
```

```
35    y.max         (dens$y) *
      <-      max1.1
36
37    x.to.plot.true <- x.to.plot <-
      dens$x
38
      y.seq      (0, y.max/2,
39    <-      cy.max)
40    y.to.plot <-
      dens$y
41
      range.x
42    =         range(x)
43    cum.dens
      <-        ecdf(x)(x.to.plot)
44

45

46

      n.blanks <-
47    10
48    cum.dens           , n.blanks),
49    <-        c(rep(NAcum.dens)
50    diff <- x.to.plot[1] -
      x.to.plot[2]
51
      blanks <- diff * (n.blanks:1) +
52    x.to.plot[1]
53    x.to.plot     (blanks,
      <-        cx.to.plot)
54
      y.to.plot        (0, n.blanks),
55    <-        c(repy.to.plot)
56

57
      cum.perc <- cum.dens *
58    100
59    z.mat          (cum.perc, byrow      , nrow = 3, ncol = n +
      <-      matrix=                TRUEn.blanks,
60
                     dimnames       (y = y.seq, x =
61                   =        listx.to.plot))
62

63    col.fun <-                      (colors, domain = 0:1, na.color
      scales::          col_numeric=                          "white")
64
      x.as.colors
65    <-        col.fun(cum.dens)
66
```

```
67    z.to.plot.scaled <-
      scales::                      rescale(cum.perc)
68
      color.lookup                        (z.to.plot.scaled,
69    <-            setNames(data.framex.as.colors),            NULL)
70

71    require(plotly)
72
      p              (z =
73    <-  plot_lyz.mat,
74                 xsrc =
                   x.to.plot,
75
                   ysrc =
76                 y.seq,
77                 type
                   =     "heatmap",
78
                   colorscale =
79                 color.lookup,
80                 cauto
                   =      FALSE,
81
                   hoverinfo
82                 =           "none",
83
                    colorbar       (title =
84                 =           listlegend.title),
85                 showscale =
                   show.legend)
86

87
      p
88    <-  add_trace(p,
89
                   x    (1:(n + n.blanks), (n +
                   =  cn.blanks):1),
                   y                  (y.max * 1.10, n +
                   =  c(y.to.plot,repn.blanks)),
                   fill
                   =     "tonexty",
                   hoverinfo
                   =           "none",
                   showlegend
                   =           FALSE,
                   type
                   =     "scatter",
```

```r
                mode
                =        "line",

                showscale
                =                FALSE,

                line        (color              , width =
                =      list=          "white"0),

                fillcolor
                =                "white")


    p
    <-   add_trace(p,

                x = 1:(n +
                n.blanks),

                y =
                y.to.plot,

                name
                =        "",

                hoverinfo
                =                "text",

                text                              ": %.0f %% <
                =      sprintf(paste0(x.title,%.1f"
), cum.perc,
x.to.plot),

                type
                =        "scatter",

                mode
                =        "lines",

                line        (color            , width =
                =      list=          "white"0),

                showlegend=FALSE,

                showscale=FALSE)

    p <-                (p, displayModeBar
    plotly::      config=                        FALSE)


    x.text          (x.to.plot, n =
    <-        prettyn.breaks)


x.tick <- 1 + (x.text - x.to.plot[1]) / (x.to.plot[n + n.blanks] -
x.to.plot[1]) * (n + n.blanks - 1)
```

```
    p              (p, title =
    <-   layouttitle,

              xaxis       (title = x.title, tickmode
              =      list=                         "array"
, tickvals = x.tick, ticktext =
x.text),

              yaxis       (title    , showline       , ticks
              =      list=       ""=          FALSE=        ""
, showticklabels      ,           (0,
=              FALSErange=  cy.max)),

              margin      (t = 30, l = 5, b = 50, r =
              =       list5))

    p

}
```

---

## Acknowledgements

My colleague Carmen wrote most of the code. It is written using the wonderful plotly package.

### Author: Tim Bock

Tim Bock is the founder of Displayr. Tim is a data scientist, who has consulted, published academic papers, and won awards, for problems/techniques as diverse as neural networks, mixture models, data fusion, market segmentation, IPO pricing, small sample research, and data visualization. He has conducted data science projects for numerous companies, including Pfizer, Coca Cola, ACNielsen, KFC, Weight Watchers, Unilever, and Nestle. He is also the founder of Q www.qresearchsoftware.com, a data science product designed for survey research, which is used by all the world's seven largest market research consultancies. He studied econometrics, maths, and marketing, and has a University Medal and PhD from the University of New South Wales (Australia's leading research university), where he was an adjunct member of staff for 15 years.