# Buzzfeed Data Pairs Matrix Code

By Max Woolf (http://minimaxir.com)

This notebook is the complement to my blog post Facebook Reactions and the Problem With Quantifying Likes Differently.

*This notebook is licensed under the MIT License. If you use the code or data visualization designs contained within this notebook, it would be greatly appreciated if proper attribution is given back to this notebook and/or myself. Thanks! :)*

```r
1  options(warn = -1)
2
3  # IMPORTANT: This assumes that all packages in "Rstart.R" are installed,
4  # and the fonts "Source Sans Pro" and "Open Sans Condensed Bold" are installed
5  # via extrafont. If ggplot2 charts fail to render, you may need to change/remove the theme
       call.
6
7  source("Rstart.R")
8  library(GGally) # ggpairs
9
10 sessionInfo()
```

```
1  Attaching package: ''dplyr
2
3  The following objects are masked from 'package:'stats:
4
5      filter, lag
6
7  The following objects are masked from 'package:'base:
8
9      intersect, setdiff, setequal, union
10
11 Registering fonts with R
12
13 Attaching package: ''scales
14
15 The following objects are masked from 'package:'readr:
16
17     col_factor, col_numeric
18
19
20 Attaching package: ''GGally
21
22 The following object is masked from 'package:'dplyr:
23
24     nasa
25
26
27
28
29
30
31 R version 3.2.3 (2015-12-10)
32 Platform: x86_64-apple-darwin13.4.0 (64-bit)
33 Running under: OS X 10.11.3 (El Capitan)
```

```
34
35 locale:
36 [1] en_US.UTF-8/en_US.UTF-8/en_US.UTF-8/C/en_US.UTF-8/en_US.UTF-8
37
38 attached base packages:
39 [1] grid      stats     graphics  grDevices utils     datasets  methods
40 [8] base
41
42 other attached packages:
43 [1] GGally_1.0.1      stringr_1.0.0     digest_0.6.8      RColorBrewer_1.1-2
44 [5] scales_0.3.0      extrafont_0.17    ggplot2_2.0.0     dplyr_0.4.3
45 [9] readr_0.1.1
46
47 loaded via a namespace (and not attached):
48  [1] Rcpp_0.12.1      Rttf2pt1_1.3.3   magrittr_1.5     munsell_0.4.2
49  [5] uuid_0.1-2       colorspace_1.2-6 R6_2.1.1         plyr_1.8.3
50  [9] tools_3.2.3      parallel_3.2.3   gtable_0.1.2     DBI_0.3.1
51 [13] extrafontdb_1.0  assertthat_0.1   IRdisplay_0.3    repr_0.4
52 [17] base64enc_0.1-3  IRkernel_0.5     evaluate_0.8     rzmq_0.7.7
53 [21] stringi_0.5-5    reshape_0.8.5    jsonlite_0.9.19
```

```r
1 df <- read_csv("buzzfeed_data_social_10k.csv")
2
3 print(df)
```

```
1 Source: local data frame [10,388 x 22]
2
3                                                             title
4                                                             (chr)
5 1                        How Well Do You Know Your Banned Books?
6 2   16 Things F. Scott Fitzgerald Doesn't Want You To Worry About
7 3            Watch Nick And Amy's Fatal Attraction In "Gone Girl"
8 4   Alison Bechdel Is The Ultimate Genius "Dyke To Watch Out For"
9 5                         16 Reasons You'd Probably Die At Hogwarts
10 6              19 Banned Books If They Were Made Appropriate
11 7                            "Zelda's Dreams," By James Franco
12 8                        How Scandalous Is Your Reading History?
13 9                "Gone Girl" Is Now A Sleek But Hollow Movie
14 10           17 Things English Majors Are Tired Of Hearing
15 ..                                                            ...
16 Variables not shown: url (chr), author (chr), date (date), category (chr),
17   special (chr), responses (int), num_fb_shares (int), num_tweets (int),
18   num_fb_comments (int), love (int), yaaass (int), helpful (int), omg (int),
19   lol (int), cute (int), win (int), wtf (int), fail (int), trashy (int), ew
20   (int), hate (int)
```

Select only the columns with reaction data, and get spot correlations.

```r
1 df_reactions <- na.omit(df %>% select(love:hate))
2
3 print(df_reactions)
4
5 print(cor(df_reactions))
```

```
1 Source: local data frame [9,883 x 12]
```

```
 2
 3      love yaaass helpful    omg   lol  cute   win   wtf  fail trashy    ew  hate
 4     (int)  (int)   (int)  (int) (int) (int) (int) (int) (int)  (int) (int) (int)
 5 1     31      0       3      7     1     1     3     5     4      0     0     1
 6 2    110      0       0      2     9    17    18     7     0      1     0     0
 7 3      5      0       0      0     0     0     2     0     0      0     0     0
 8 4     16      0       0      0     0     0     1     0     0      0     0     0
 9 5     72      0       0      2    25     1     4     0     4      0     0     0
10 6     44      7       0      4    20     1     8     3     7      1     0     0
11 7     25      0       0      0     0     0     0     7     2      0     0     0
12 8    139      2       1      5    10     1    20     1     0      2     0     1
13 9     19      0       0      2     2     0     1     0     0      0     0     0
14 10   119     23       2      3    22     1    25     0     1      0     0     0
15 ..    ...    ...     ...    ...   ...   ...   ...   ...   ...    ...   ...   ...
16               love      yaaass      helpful         omg         lol        cute
17 love    1.00000000 0.46626799 0.124755232 0.68036925 0.47360895  0.629094452
18 yaaass  0.46626799 1.00000000 0.175511580 0.35403737 0.26705946  0.096387912
19 helpful 0.12475523 0.17551158 1.000000000 0.04352179 0.01926325  0.008081270
20 omg     0.68036925 0.35403737 0.043521787 1.00000000 0.38471634  0.539838706
21 lol     0.47360895 0.26705946 0.019263247 0.38471634 1.00000000  0.305064425
22 cute    0.62909445 0.09638791 0.008081270 0.53983871 0.30506443  1.000000000
23 win     0.83126618 0.45288311 0.114922581 0.59319268 0.43868351  0.523278287
24 wtf     0.09907593 0.05272187 0.022267654 0.31346725 0.20007750  0.008643063
25 fail    0.07005368 0.07472599 0.021095192 0.18130431 0.17963674 -0.031200420
26 trashy  0.03739368 0.09077292 0.014492817 0.13685420 0.09558570 -0.031452816
27 ew      0.05038921 0.10098157 0.009602044 0.20792642 0.11147785 -0.024514311
28 hate    0.15831206 0.02737651 0.015482722 0.27294572 0.05448569  0.007341387
29               win         wtf        fail      trashy          ew        hate
30 love    0.83126618 0.099075927  0.07005368  0.03739368  0.050389209 0.158312061
31 yaaass  0.45288311 0.052721871  0.07472599  0.09077292  0.100981567 0.027376513
32 helpful 0.11492258 0.022267654  0.02109519  0.01449282  0.009602044 0.015482722
33 omg     0.59319268 0.313467249  0.18130431  0.13685420  0.207926422 0.272945720
34 lol     0.43868351 0.200077499  0.17963674  0.09558570  0.111477851 0.054485686
35 cute    0.52327829 0.008643063 -0.03120042 -0.03145282 -0.024514311 0.007341387
36 win     1.00000000 0.061382292  0.04877020  0.02347292  0.023725465 0.070561338
37 wtf     0.06138229 1.000000000  0.63592405  0.50851441  0.566388147 0.332060843
38 fail    0.04877020 0.635924055  1.00000000  0.51560199  0.505881072 0.348757439
39 trashy  0.02347292 0.508514410  0.51560199  1.00000000  0.805459962 0.255968387
40 ew      0.02372546 0.566388147  0.50588107  0.80545996  1.000000000 0.255072265
41 hate    0.07056134 0.332060843  0.34875744  0.25596839  0.255072265 1.000000000
```

Note that the `helpful` and `trashy` reactions are not used in 2016, so we will not use them.

Use `ggpairs` to plot multidimensional data (lower and diag functions adapted from the GGally package viginette; upper correlation function adopted from Barret Schloerke on GitHub).

```
1 pairs_theme <- function (x) {
2                 theme_bw(base_size = 5) +
3                 theme(panel.grid.minor.x = element_blank()) +
4                 theme(panel.grid.minor.y = element_blank())
5                 }
6
7
8 gglower <- function(data, mapping, ..., high = "#c0392b") {
9   ggplot(data = data, mapping = mapping) +
```

```
10      geom_bin2d(...) +
11      scale_x_log10(limits=c(10^0,10^3), breaks=10^(0:3)) +
12      scale_y_log10(limits=c(10^0,10^3), breaks=10^(0:3)) +
13      geom_smooth(alpha = 0.5, size = 0.25, color = "#1a1a1a", method = "lm") +
14      scale_fill_gradient(low = "#EEEEEE", high = high, trans = "log") +
15      pairs_theme()
16 }
17
18 ggdiag <- function(data, mapping, ..., color = "#1a1a1a") {
19   ggplot(data = data, mapping = mapping) +
20      geom_density(..., color = color) +
21      scale_x_log10(limits=c(10^0,10^3), breaks=10^(0:3)) +
22      pairs_theme()
23 }
24
25 # From https://github.com/ggobi/ggally/issues/139#issuecomment-176271618
26
27 ggupper <- function(data, mapping, color = I("grey50"), sizeRange = c(1, 3), ...) {
28
29   # get the x and y data to use the other code
30   x <- eval(mapping$x, data)
31   y <- eval(mapping$y, data)
32
33   ct <- cor.test(x,y)
34   sig <- symnum(
35     ct$p.value, corr = FALSE, na = FALSE,
36     cutpoints = c(0, 0.001, 0.01, 0.05, 0.1, 1),
37     symbols = c("***", "**", "*", ".", " ")
38   )
39
40   r <- unname(ct$estimate)
41   rt <- format(r, digits=2)[1]
42
43   # since we can't print it to get the strsize, just use the max size range
44   cex <- max(sizeRange)
45
46   # helper function to calculate a useable size
47   percent_of_range <- function(percent, range) {
48     percent * diff(range) + min(range, na.rm = TRUE)
49   }
50
51   # plot the cor value
52   ggally_text(
53     label = as.character(rt),
54     mapping = aes(),
55     xP = 0.5, yP = 0.5,
56     size = I(percent_of_range(cex * abs(r), sizeRange)),
57     color = color,
58     ...
59   ) +
60     # add the sig stars
61     geom_text(
62       aes_string(
63         x = 0.8,
```

```
64        y = 0.8
65      ),
66      label = sig,
67      size = I(cex),
68      color = color,
69      ...
70    ) +
71    pairs_theme() +
72    theme(panel.grid.major.x = element_blank()) +
73    theme(panel.grid.major.y = element_blank())
74
75 }
```

```
1 pos_color <- "#27ae60"
2
3 plot <- ggpairs(df_reactions, columns = c("love", "yaaass", "omg", "lol", "cute", "win"),
4        title = sprintf("Pairs Plot of Positive Reaction Counts on %00d BuzzFeed Articles",
              nrow(df_reactions)),
5        upper = list(continuous = wrap(ggupper, color = pos_color)),
6        lower = list(continuous = wrap(gglower, high = pos_color)),
7        diag = list(continuous = wrap(ggdiag, color = pos_color))) +
8        theme(title = element_text(size=10))
9
10 png("buzzfeed-pos.png", w=1600, h=1600, res=300)
11 plot
12 dev.off()
```

pdf: 2

```
1 neg_color <- "#c0392b"
2
3 plot <- ggpairs(df_reactions, columns = c("love", "wtf", "fail", "ew", "hate"),
4        title = sprintf("Pairs Plot of Love + Negative Reaction Counts on %00d BuzzFeed
              Articles", nrow(df_reactions)),
5        upper = list(continuous = wrap(ggupper, color = neg_color )),
6        lower = list(continuous = wrap(gglower, high = neg_color)),
7        diag = list(continuous = wrap(ggdiag, color = neg_color))) +
8        theme(title = element_text(size=10))
9
10 png("buzzfeed-neg.png", w=1600, h=1600, res=300)
11 plot
12 dev.off()
```
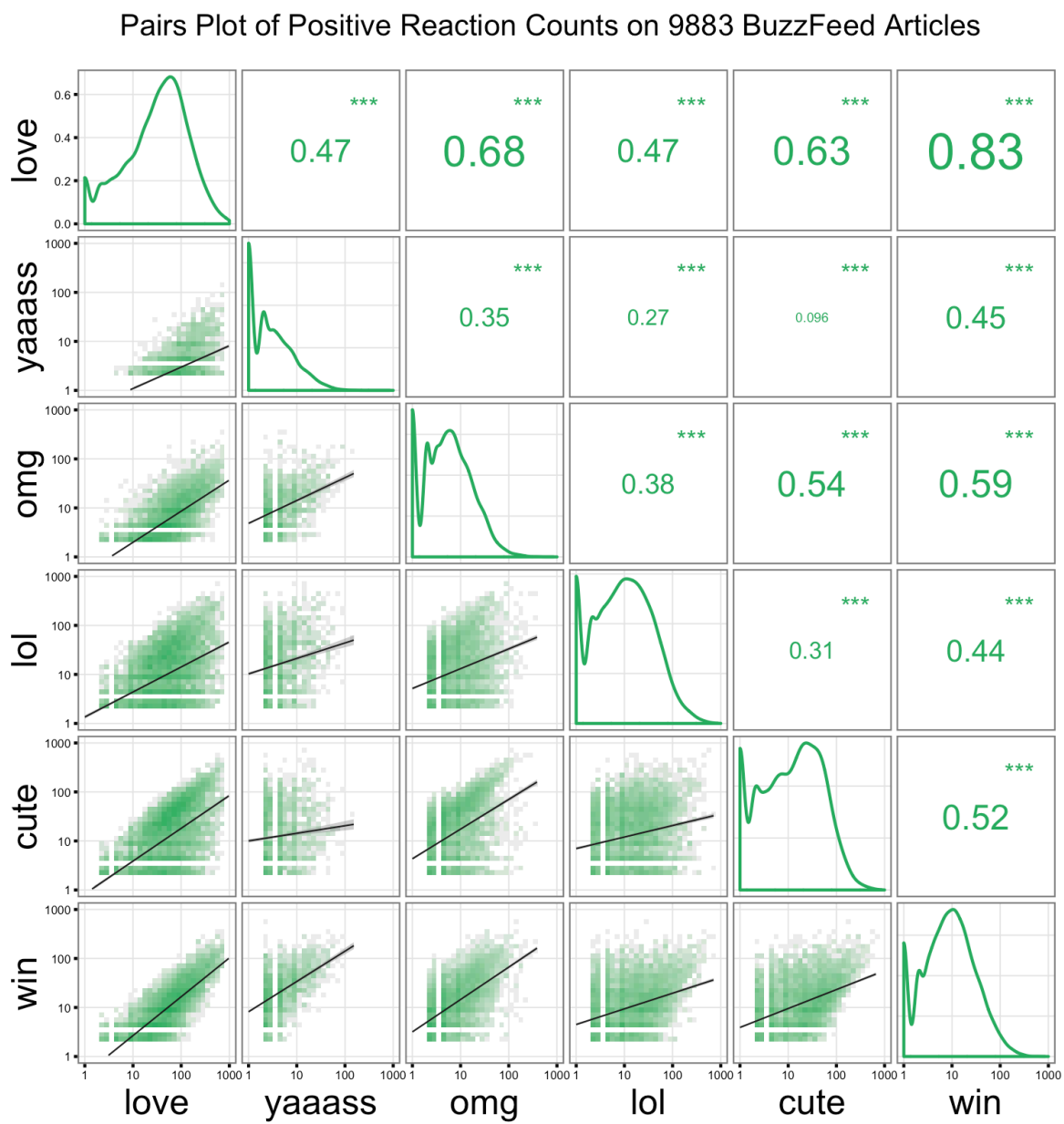
pdf: 2

# The MIT License (MIT)

Figure 1:

Pairs Plot of Love + Negative Reaction Counts on 9883 BuzzFeed Articles
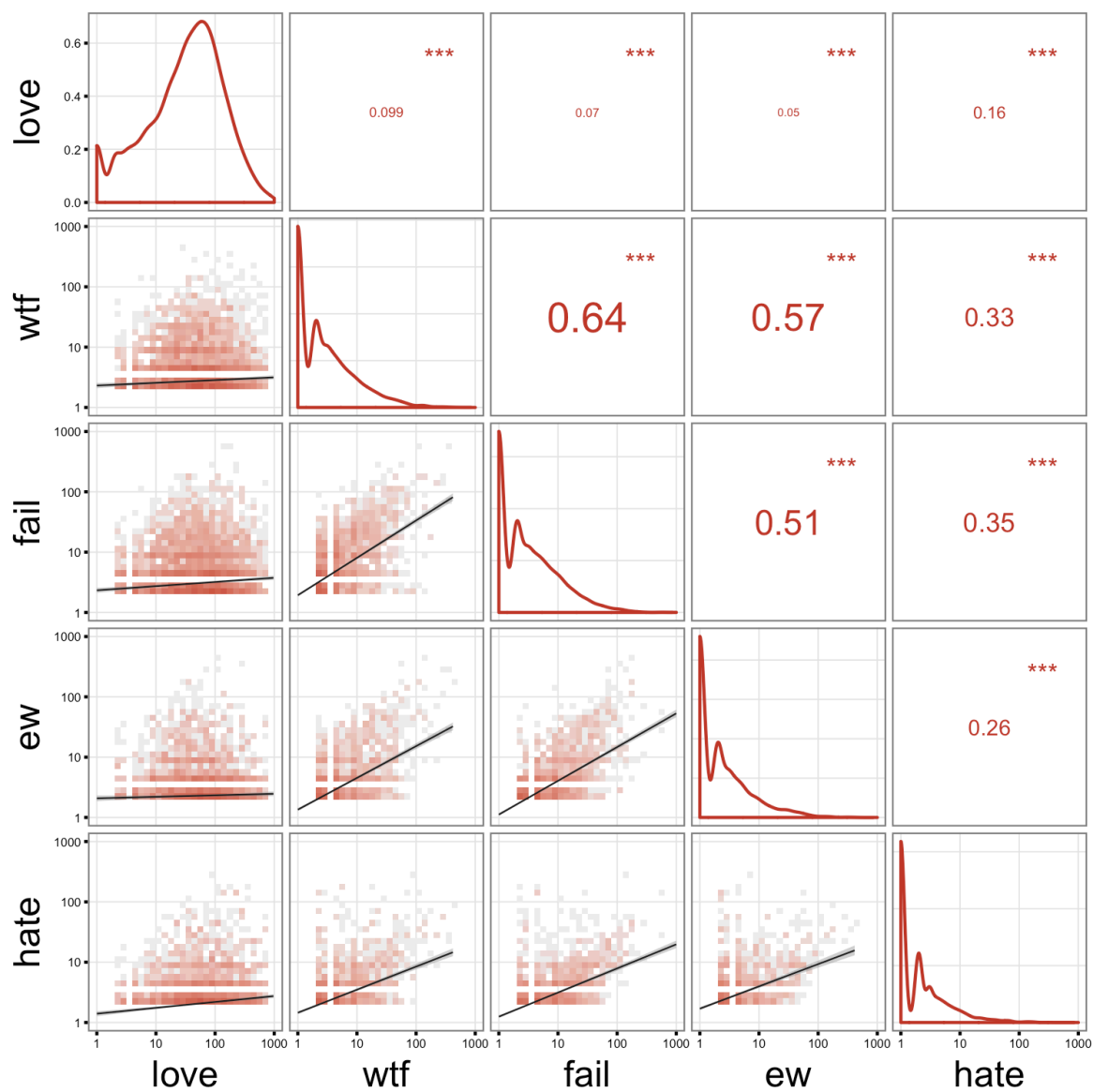
Figure 2:

The above copyright notice and this permission notice shall be included in all copies or substantial portions of the Software.

THE SOFTWARE IS PROVIDED "AS IS", WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT. IN NO EVENT SHALL THE AUTHORS OR COPYRIGHT HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR OTHER DEALINGS IN THE SOFTWARE.