# Spotify Recommendations

*Duncan Gates*

*30 November, 2020*

## Some quick cleaning

The dataset came from Kaggle but was still a little messy, not sure why the artists column came wrapped in ['name']. Additionally some date formatting was inconsistent (as always), I imputed -01-01 if there was no date after the year.
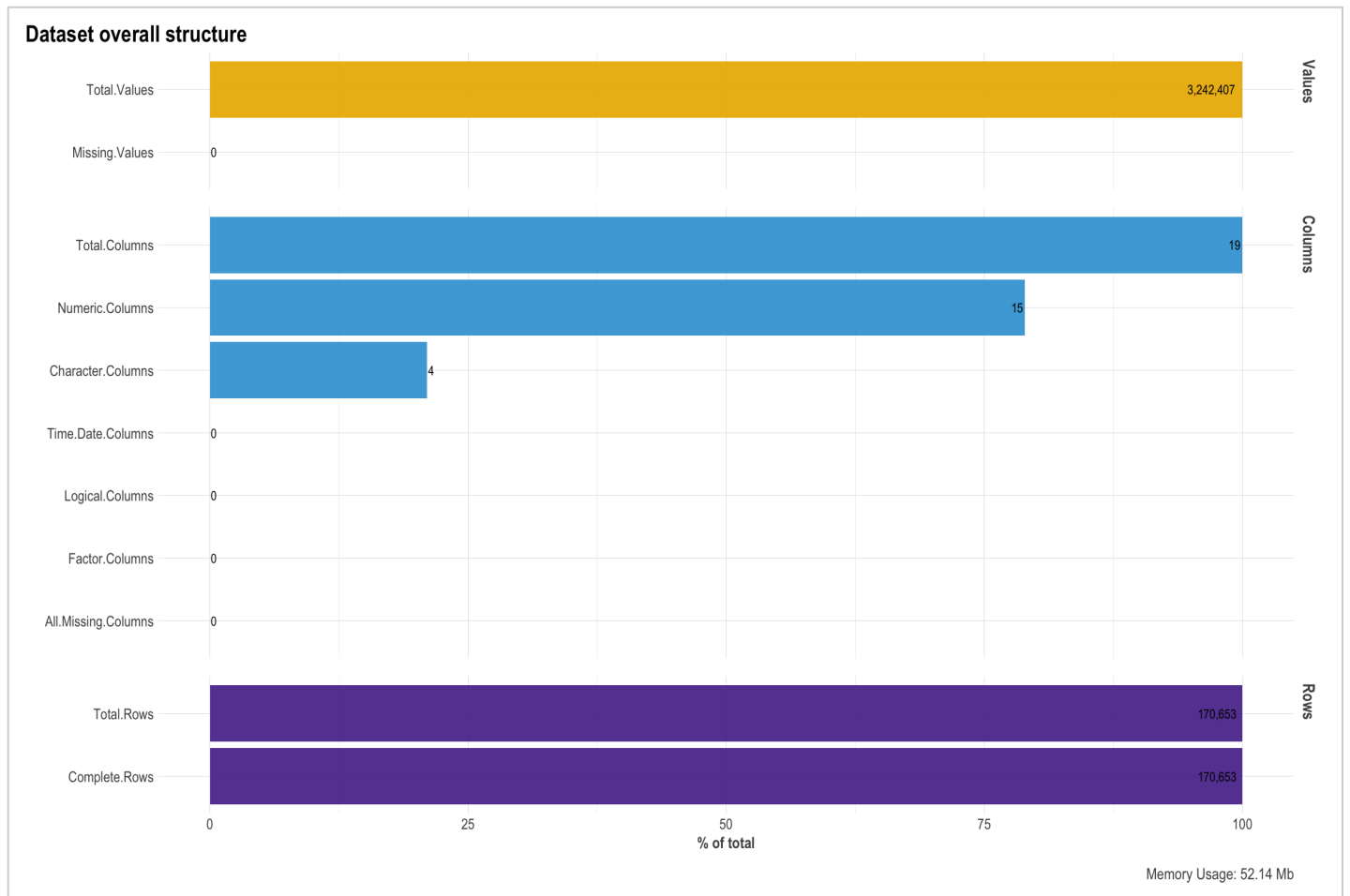
```
df2 <- df %>% mutate(artists = str_remove_all(artists, "\\['"),
                         artists = str_remove_all(artists, "\\']"),
                           artists = gsub(",", " and ", artists),
                  artists = str_remove_all(artists, "^'|'$"), # work on this a bit mo
                                          re
                     decade = as.factor(floor(year/10)*10), # make decade column
                                  year = as.factor(year),
                     release_date = ifelse(nchar(release_date) == 4, paste0(release_dat
            e, "-01-01"), release_date)) # Make release_date a real date column
      # Also would be cool to try to recognize gender by name and make a dummy column
```

This one is for networking later.

```
df_genre2 <- df_genre %>% mutate(genres = str_remove_all(genres, "\\["), # Get rid of br
                                        ackets
                            genres = str_remove_all(genres, "\\]"),
                  genres = str_split(genres, ",")) %>%  # Split up genres by comma
                                  unnest(genres) %>%
               mutate(genres = str_remove_all(genres, "\\ '")) %>% # remove all sp
                               ace appostrophes
                   mutate(genres = str_remove_all(genres, "\\'")) %>% #remove a
             ll appostrophes, could be a better filter than this
                 mutate(genres = str_remove_all(genres, "\"")) # Some quotes stuck a
                     round on childrens music category
    # df_genre3 <- df_genre2 %>% pivot_wider(names_from = artists, values_from = genres)
```

## Looking At the Data Using Lares

```
df_str(df, return = "plot")
```

**Dataset overall structure**



Memory Usage: 52.14 Mb

Look at percentages and cumulatives, looks like no popularity is very common!

```
df %>% freqs(popularity, plot = T)
```

**Frequencies and Percentages**

*Variable: popularity [20 most frequent]*

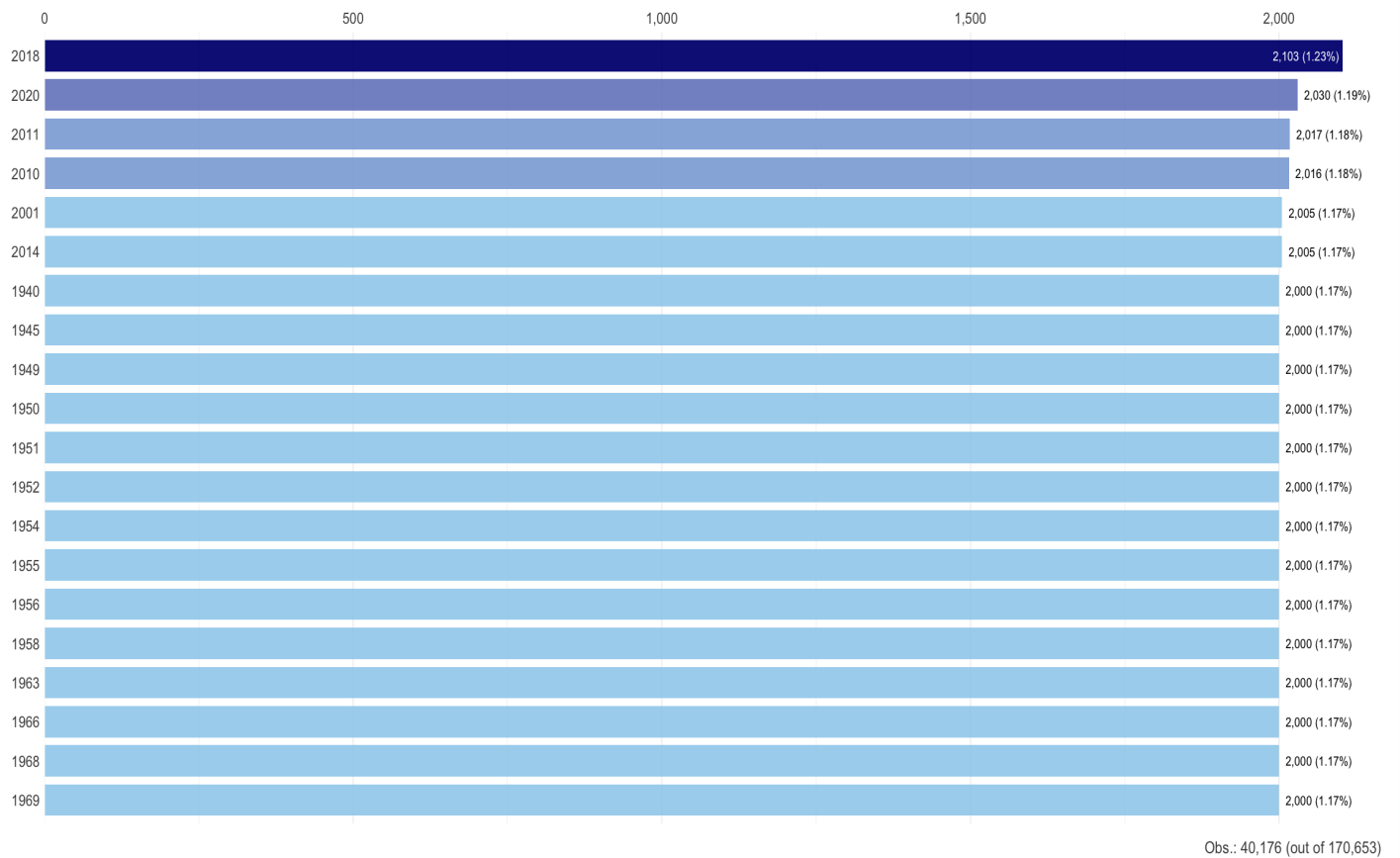| | | |
|---|---|---|
| 0 | | 27,892 (16.34%) |
| 43 | | 3,136 (1.84%) |
| 44 | | 3,117 (1.83%) |
| 41 | | 3,078 (1.8%) |
| 40 | | 3,051 (1.79%) |
| 42 | | 3,051 (1.79%) |
| 39 | | 2,968 (1.74%) |
| 36 | | 2,896 (1.7%) |
| 46 | | 2,883 (1.69%) |
| 1 | | 2,876 (1.69%) |
| 35 | | 2,864 (1.68%) |
| 37 | | 2,836 (1.66%) |
| 45 | | 2,832 (1.66%) |
| 47 | | 2,806 (1.64%) |
| 34 | | 2,803 (1.64%) |
| 38 | | 2,749 (1.61%) |
| 48 | | 2,726 (1.6%) |
| 49 | | 2,689 (1.58%) |
| 32 | | 2,625 (1.54%) |
| 31 | | 2,615 (1.53%) |

Obs.: 82,493 (out of 170,653)

There is about the same amount of music every year in this dataframe

```
df %>% freqs(year, plot = T)
```

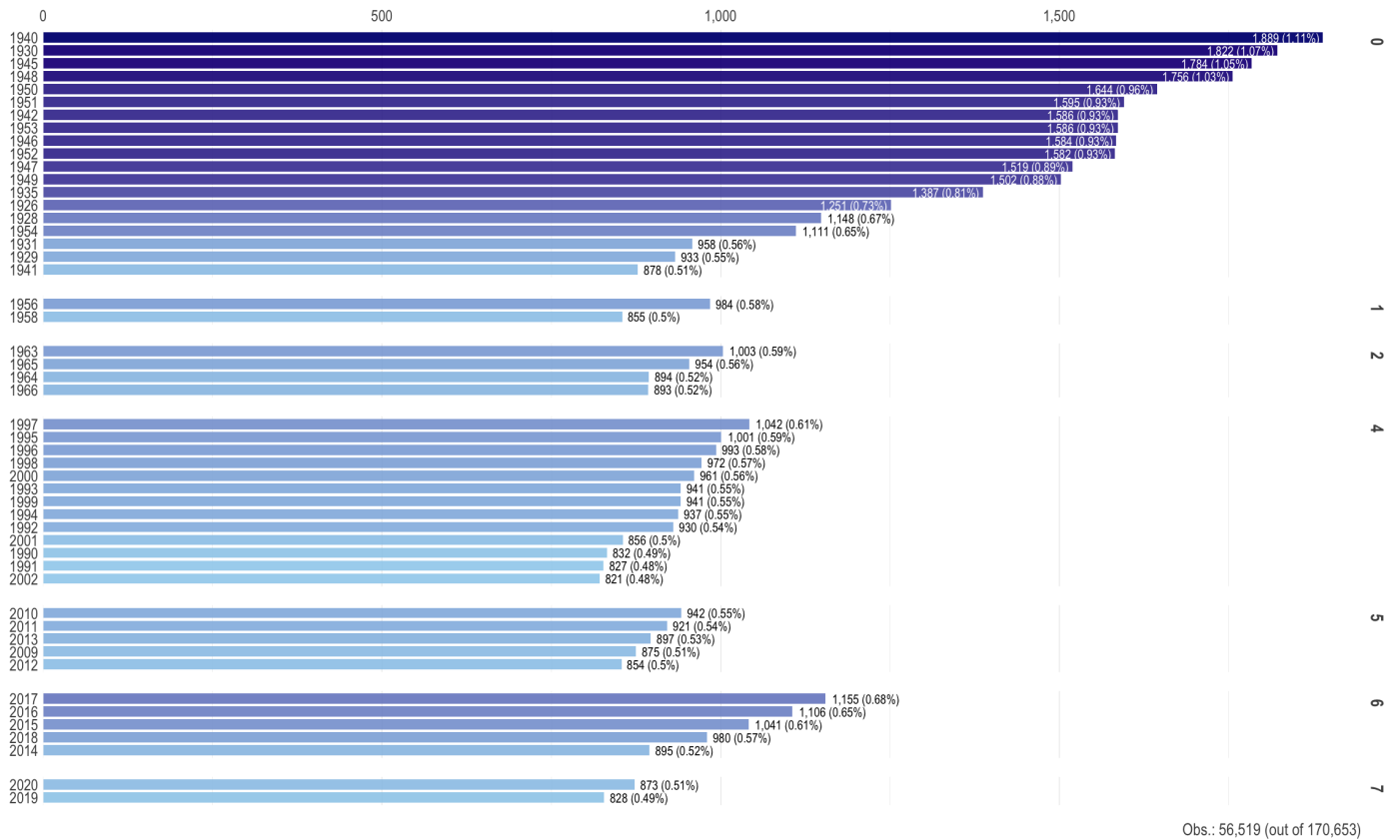## Frequencies and Percentages
*Variable: year [20 most frequent]*

| Year | | Frequency |
|------|---|-----------|
| 2018 | | 2,103 (1.23%) |
| 2020 | | 2,030 (1.19%) |
| 2011 | | 2,017 (1.18%) |
| 2010 | | 2,016 (1.18%) |
| 2001 | | 2,005 (1.17%) |
| 2014 | | 2,005 (1.17%) |
| 1940 | | 2,000 (1.17%) |
| 1945 | | 2,000 (1.17%) |
| 1949 | | 2,000 (1.17%) |
| 1950 | | 2,000 (1.17%) |
| 1951 | | 2,000 (1.17%) |
| 1952 | | 2,000 (1.17%) |
| 1954 | | 2,000 (1.17%) |
| 1955 | | 2,000 (1.17%) |
| 1956 | | 2,000 (1.17%) |
| 1958 | | 2,000 (1.17%) |
| 1963 | | 2,000 (1.17%) |
| 1966 | | 2,000 (1.17%) |
| 1968 | | 2,000 (1.17%) |
| 1969 | | 2,000 (1.17%) |

Obs.: 40,176 (out of 170,653)

Basically newer stuff is more popular with little to no exceptions

```
df %>% mutate(popularity = round(popularity/10)) %>% freqs(popularity, year, plot = T, r
esults = F, top = 50)
```
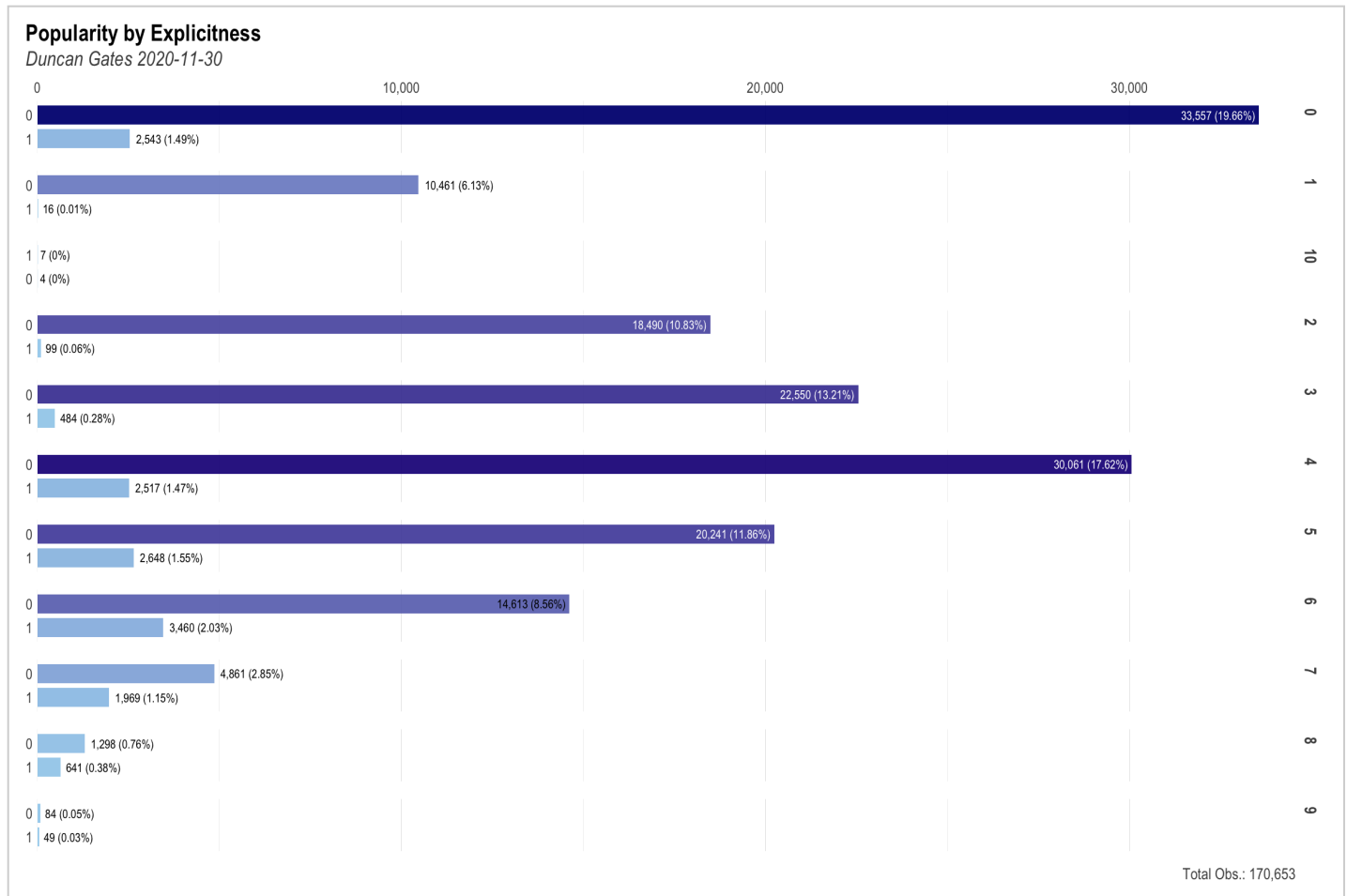
**Frequencies and Percentages**
*Variables: year grouped by popularity*
*[50 most frequent]*



| | 0 | 500 | 1,000 | 1,500 | |
|---|---|---|---|---|---|

```
1940                                                    1,889 (1.11%)
1930                                                  1,822 (1.07%)
1945                                                 1,784 (1.05%)
1948                                                1,756 (1.03%)      0
1950                                            1,644 (0.96%)
1951                                         1,595 (0.93%)
1942                                         1,586 (0.93%)
1953                                         1,586 (0.93%)
1946                                         1,584 (0.93%)
1952                                         1,582 (0.93%)
1947                                       1,519 (0.89%)
1949                                      1,502 (0.88%)
1935                                 1,387 (0.81%)
1926                          1,251 (0.73%)
1928                      1,148 (0.67%)
1954                     1,111 (0.65%)
1931               958 (0.56%)
1929              933 (0.55%)
1941            878 (0.51%)

1956               984 (0.58%)                                          1
1958           855 (0.5%)

1963               1,003 (0.59%)                                        2
1965              954 (0.56%)
1964            894 (0.52%)
1966            893 (0.52%)

1997                1,042 (0.61%)                                       4
1995               1,001 (0.59%)
1996              993 (0.58%)
1998              972 (0.57%)
2000             961 (0.56%)
1993            941 (0.55%)
1999            941 (0.55%)
1994            937 (0.55%)
1992            930 (0.54%)
2001          856 (0.5%)
1990         832 (0.49%)
1991         827 (0.48%)
2002        821 (0.48%)

2010            942 (0.55%)                                             5
2011           921 (0.54%)
2013           897 (0.53%)
2009          875 (0.51%)
2012          854 (0.5%)

2017                1,155 (0.68%)                                       6
2016               1,106 (0.65%)
2015              1,041 (0.61%)
2018             980 (0.57%)
2014           895 (0.52%)

2020          873 (0.51%)                                               7
2019         828 (0.49%)
```

Obs.: 56,519 (out of 170,653)

Looks like explicitness has a normal distribution compared to popularity

```r
df %>%
  mutate(explicit = as.factor(explicit),
    popularity = round(popularity/10)) %>%
  freqs(popularity, explicit, plot = T,
    title = "Popularity by Explicitness",
  subtitle = paste("Duncan Gates", Sys.Date()),
            results = F)
```
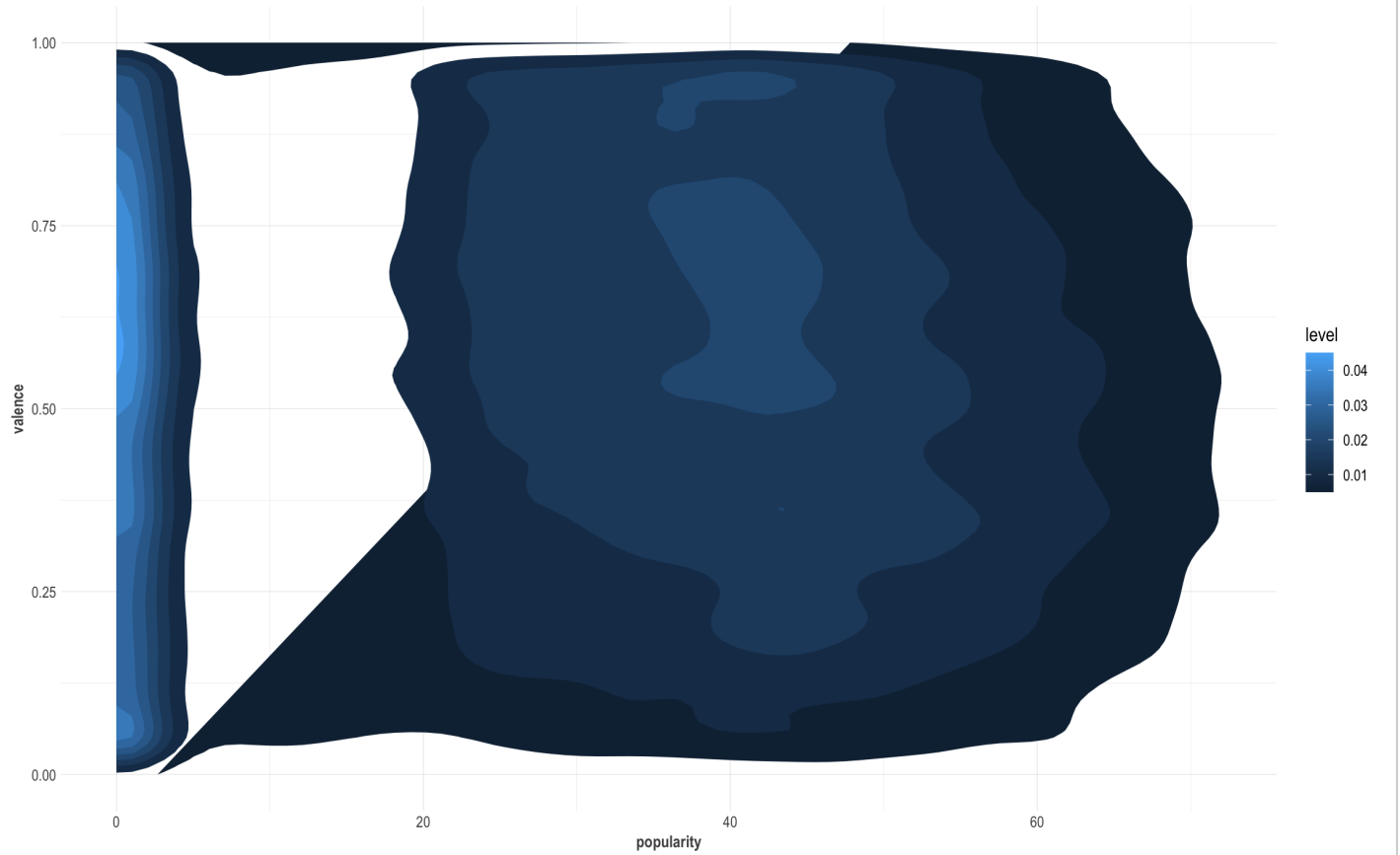
**Popularity by Explicitness**
*Duncan Gates 2020-11-30*

| | 0 | 10,000 | 20,000 | 30,000 |
|---|---|---|---|---|

**0**
0 — 33,557 (19.66%)
1 — 2,543 (1.49%)

**1**
0 — 10,461 (6.13%)
1 — 16 (0.01%)

**10**
1 — 7 (0%)
0 — 4 (0%)

**2**
0 — 18,490 (10.83%)
1 — 99 (0.06%)

**3**
0 — 22,550 (13.21%)
1 — 484 (0.28%)

**4**
0 — 30,061 (17.62%)
1 — 2,517 (1.47%)

**5**
0 — 20,241 (11.86%)
1 — 2,648 (1.55%)

**6**
0 — 14,613 (8.56%)
1 — 3,460 (2.03%)

**7**
0 — 4,861 (2.85%)
1 — 1,969 (1.15%)

**8**
0 — 1,298 (0.76%)
1 — 641 (0.38%)

**9**
0 — 84 (0.05%)
1 — 49 (0.03%)

Total Obs.: 170,653

Now we check out the distribution, there's some really cool stuff here

```
df %>% distr(popularity, valence) # Some really cool density plots can be done here
```

**2D Density Distribution**

*Variables: valence vs. popularity. Obs: 170,653*



There's also some really long songs out there…

```
df %>% distr(duration_ms)
```

**Density Distribution**
*Variable: duration_ms*

0.000015

0.000010

0.000005

0.000000

0          2000000          4000000

Obs: 170,653

This looks more like the actual distribution

```
df %>%
  filter(duration_ms < 400000) %>%
    distr(duration_ms)
```

**Density Distribution**
*Variable: duration_ms*

Obs: 161,250

Very interesting distribution here

```
df %>% distr(popularity)
```

**Density Distribution**
*Variable: popularity*



Obs: 170,653

Looks like things are a lot more explicit in 2000-2020 as one might expect, would be interesting to see how when this starts, or what drives it. I also wonder what happened in 1920-1940?

```
df %>%
  mutate(explicit = as.factor(explicit),
  new_era = ifelse(year %in% c(2000:2020), 1, 0)) %>%
    distr(explicit, new_era)
```

**Distribution and Proportions**

*Variables: explicit vs. new_era. Obs: 170,653*

```
df %>%
    mutate(decade = floor(year/10)*10) %>%
    distr(explicit, decade, custom_colours = T, abc = T)
```

**Distribution and Proportions**
*Variables: explicit vs. decade. Obs: 170,653*

By the way mode is just whether the song is major or minor.

```
df %>% distr(mode, force = "char")
```

**Frequencies and Percentages**
*Variable: mode*

| 0 | 25,000 | 50,000 | 75,000 | 100,000 |
|---|--------|--------|--------|---------|

1                                                                120,635 (70.69%)

0          50,018 (29.31%)

Total Obs.: 170,653

You can even use ggplot2!

```
                            df %>%
        distr(popularity, loudness) + geom_point(color = "yellow")
```

**2D Density Distribution**
*Variables: loudness vs. popularity. Obs: 170,653*

Wouldn't be data science without some random regressions, even more data science/machine learningy since the second one is a log odds table!

```
df %>%
  select(-c(id, name, artists, year, release_date, key)) %>%
  corr_cross(top = 10) # Look at top 10 correlations in the data, key messes with this i
                        dk why
```

## Ranked Cross-Correlations
*10 most relevant*



| variables | | |
|---|---|---|
| energy + loudness | | 78.2 |
| acousticness + energy | | -74.9 |
| acousticness + popularity | | -57.3 |
| acousticness + loudness | | -56.2 |
| valence + danceability | | 55.9 |
| energy + popularity | | 48.5 |
| loudness + popularity | | 45.7 |
| explicit + speechiness | | 41.4 |
| instrumentalness + loudness | | -40.9 |
| valence + energy | | 35.4 |

Correlation [%]

```
table <- df %>%
  select(-c(id, name, artists, year, release_date, key)) %>%
  corr_var(popularity, logs = T, plot = F, top = 10)
table %>%
mutate(variables = paste(toupper(substr(variables, 1, 1)), substr(variables, 2, nchar(
  as.character(variables))), sep = "")) %>%
mutate(corr = kableExtra::cell_spec(corr, "html", color = ifelse(corr > 0, "blue", "re
  d")),
  pvalue = kableExtra::cell_spec(pvalue, "html", color = ifelse(pvalue < 0.05, "g
  reen", "red"))) %>%
  kableExtra::kable(format = "html", escape = F) %>%
  kableExtra::kable_styling("striped", full_width = F, position = "center")
```

| variables | corr | pvalue |
|---|---|---|
| Popularity_log | 0.890732 | 0 |
| Acousticness | -0.573162 | 0 |
| Acousticness_log | -0.55757 | 0 |
| Energy_log | 0.488822 | 0 |
| Energy | 0.485005 | 0 |
| Loudness | 0.457051 | 0 |
| Instrumentalness_log | -0.300402 | 0 |

| variables | corr | pvalue |
|---|---|---|
| Instrumentalness | -0.29675 | 0 |
| Danceability | 0.199606 | 0 |
| Danceability_log | 0.196287 | 0 |

```
# wow OHSE is pretty dope check this out with a better dataframe
```

# Fun with reactable

Using lares one more time to get an idea of the data, there's a lot of NA's at first so I drop those and look again

```
df_genre2 %>% distr(genres)
```

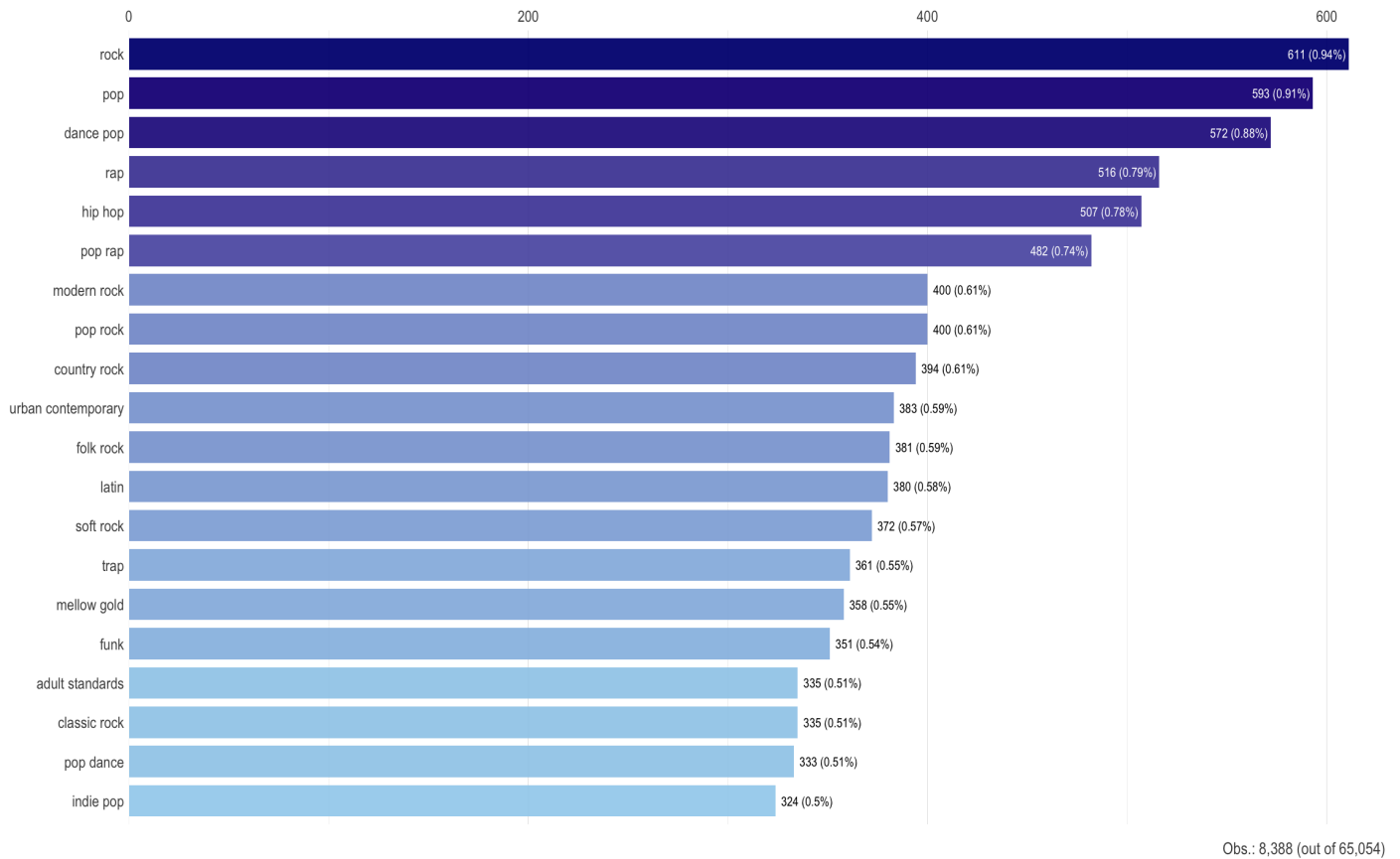**Frequencies and Percentages**
*Variable: genres [10 most frequent]*



Obs.: 14,332 (out of 74,911)

```
df_genre2 <- df_genre2 %>% na_if("") %>% na.omit %>%
                mutate(genres = as.factor(genres))
```

What's it look like for genres now? Apparently in this dataframe there has been more rock than pop, not sure if that is actually the case (it does seem possible) or if its just the nature of this data.

```
df_genre2 %>% freqs(genres, plot = T,
                    title = "Genres by Artist",
          subtitle = paste("Duncan Gates", Sys.Date()),
                    results = F)
```

**Genres by Artist**
*Duncan Gates 2020-11-30*

| Genre | Count |
|-------|-------|
| rock | 611 (0.94%) |
| pop | 593 (0.91%) |
| dance pop | 572 (0.88%) |
| rap | 516 (0.79%) |
| hip hop | 507 (0.78%) |
| pop rap | 482 (0.74%) |
| modern rock | 400 (0.61%) |
| pop rock | 400 (0.61%) |
| country rock | 394 (0.61%) |
| urban contemporary | 383 (0.59%) |
| folk rock | 381 (0.59%) |
| latin | 380 (0.58%) |
| soft rock | 372 (0.57%) |
| trap | 361 (0.55%) |
| mellow gold | 358 (0.55%) |
| funk | 351 (0.54%) |
| adult standards | 335 (0.51%) |
| classic rock | 335 (0.51%) |
| pop dance | 333 (0.51%) |
| indie pop | 324 (0.5%) |

Obs.: 8,388 (out of 65,054)

```
df_genre2 %>%
  select(-count) %>%
  summer(genres) %>%
  mutate_if(is.numeric, funs(round)) %>%
  dplyr::rename(Count = n, `Duration` = duration_ms) %>%
  mutate(`Duration` = paste0(minute(seconds_to_period((`Duration`/(1000*Count)))),
                             ":",
                             dseconds(round(seconds_to_period((`Duration`/(1000*
  Count)))), digits = 2)))) %>% # Some disgusting lubridate here sorry
  rename_with(str_to_title) %>%
  mutate(Genres = str_to_title(Genres)) %>%
  mutate_at(vars(c("Acousticness", "Danceability", "Energy", "Instrumentalness", "Livene
    ss", "Loudness", "Speechiness", "Tempo", "Valence", "Popularity", "Key", "Mode"
    )), ~round((./Count), digits = 3)) %>%
  arrange(desc(Count)) %>%
  reactable(bordered = T,
            highlight = T,
            defaultColDef = colDef(align = "center",
                                   width = 150,
            footer = function(values = c("Count", "Acousticness",
  "Danceability", "Energy", "Instrumentalness", "Liveness", "Loudness", "Speechin
            ess", "Tempo", "Valence", "Popularity", "Key", "Mode")) {
            if (!is.numeric(values)) return()
  sparkline(values, type = "bar", width = 100, height = 30) # Can also do bo
                    xplots and line graphs
            }))
```

| Genres | Count | Acousticness | Danceability | Duration |
|---|---|---|---|---|
| Rock | 611 | 0.173 | 0.502 | 4:10.06s |
| Pop | 593 | 0.258 | 0.631 | 3:37.06s |
| Dance Pop | 572 | 0.182 | 0.652 | 3:49.23s |
| Rap | 516 | 0.145 | 0.727 | 4:0.34s |
| Hip Hop | 507 | 0.16 | 0.724 | 4:7.97s |
| Pop Rap | 482 | 0.156 | 0.707 | 3:55.07s |
| Modern Rock | 400 | 0.158 | 0.538 | 3:54.39s |
| Pop Rock | 400 | 0.22 | 0.545 | 3:58.57s |
| Country Rock | 394 | 0.36 | 0.553 | 3:59.78s |
| Urban Contemporary | 383 | 0.245 | 0.661 | 4:26.81s |
| | | | | |

1–10 of 2978 rows          Previous  **1**  2  3  4  5  ...  298  Next

Let's filter down a bit and make things prettier

```r
make_color_pal <- function(colors, bias = 1) {
    get_color <- colorRamp(colors, bias = bias)
    function(x) rgb(get_color(x), maxColorValue = 255)
} # Make a color function
good_color <- make_color_pal(viridis::magma(n = 12), bias = 2)
# seq(0.1, 0.9, length.out = 12) %>%
#    good_color() %>%
#    scales::show_col() # This just shows the color palette generated

color_table <- df_genre2 %>%
    select(-count) %>%
    summer(genres) %>%
    filter(n > 200) %>% # Let's get the top 60 most popular genres
    mutate_if(is.numeric, funs(round)) %>%
    dplyr::rename(Count = n, `Duration` = duration_ms) %>%
    mutate(`Duration` = paste0(minute(seconds_to_period((`Duration`/(1000*Count)))),
                                ":",
                                dseconds(round(seconds_to_period((`Duration`/(1000*
    Count))), digits = 2)))) %>% # Some disgusting lubridate here sorry
    rename_with(str_to_title) %>%
    mutate(Genres = str_to_title(Genres)) %>%
    mutate_at(vars(c("Acousticness", "Danceability", "Energy", "Instrumentalness", "Livene
        ss", "Loudness", "Speechiness", "Tempo", "Valence", "Popularity", "Key", "Mode"
        )), ~round((./Count), digits = 3)) %>%
    arrange(desc(Count))
```

```
## Grouped by: 'genres'
```

```
## Joining, by = "genres"
```

```
color_table %>%
  reactable(bordered = T,
            highlight = T,
            columns = list(
              Count = colDef(
                name = "Count",
                style = function(value) {
                  value
                  normalized <- (value - min(color_table$Count)) / (max(color_table$Coun
                    t) - min(color_table$Count))
                  color <- rev(good_color(normalized))
                  list(background = color)
                }
              )
            ),
            defaultColDef = colDef(align = "center",
                                   width = 150,
                                   footer = function(values = c("Count", "Acousticness",
    "Danceability", "Energy", "Instrumentalness", "Liveness", "Loudness", "Speechin
            ess", "Tempo", "Valence", "Popularity", "Key", "Mode")) {
                                     if (!is.numeric(values)) return()
              sparkline(values, type = "bar", width = 100, height = 30) # Can also do bo
                                     xplots and line graphs
                                   }))
```

| Genres | Count | Acousticness | Danceability | Duration |
|---|---|---|---|---|
| Rock | 611 | 0.173 | 0.502 | 4:10.06s |
| Pop | 593 | 0.258 | 0.631 | 3:37.06s |
| Dance Pop | 572 | 0.182 | 0.652 | 3:49.23s |
| Rap | 516 | 0.145 | 0.727 | 4:0.34s |
| Hip Hop | 507 | 0.16 | 0.724 | 4:7.97s |
| Pop Rap | 482 | 0.156 | 0.707 | 3:55.07s |
| Modern Rock | 400 | 0.158 | 0.538 | 3:54.39s |
| Pop Rock | 400 | 0.22 | 0.545 | 3:58.57s |
| Country Rock | 394 | 0.36 | 0.553 | 3:59.78s |
| Urban Contemporary | 383 | 0.245 | 0.661 | 4:26.81s |
|  | |  |  |  |

1–10 of 60 rows

# Machine Learning???

Let's load some networking libraries

```
library(ggraph)
library(igraph)
```

Now let's make some central nodes for our network.

```
network_df <- df_genre2 %>% select(genres, artists, popularity)
network_graph <- network_df %>%
graph_from_data_frame() # From igraph

a <- grid::arrow(type = "closed", length = unit(.10, "inches"))
# ggraph(network_graph, layout = "fr") +
#   geom_edge_link() +
#   geom_node_point(color = "lightblue", size = 5) +
#   geom_node_text(aes(label = name), vjust = 1, hjust = 1, repel = TRUE) +
#   ggtitle("Bi-Gram Network for All Songs") +
#   theme_void()
```

Some quick cleaning    ✕

Looking At the Data Using Lares

Fun with reactable

Machine Learning???