# CDC-data-analysis

## Justin Rivera

## 2023-09-30

```r
library(tidyverse)
```

```
## Warning: package 'tidyverse' was built under R version 4.1.3

## Warning: package 'tibble' was built under R version 4.1.3

## Warning: package 'tidyr' was built under R version 4.1.3

## Warning: package 'readr' was built under R version 4.1.3

## Warning: package 'purrr' was built under R version 4.1.3

## Warning: package 'dplyr' was built under R version 4.1.3

## Warning: package 'stringr' was built under R version 4.1.3

## Warning: package 'forcats' was built under R version 4.1.3

## Warning: package 'lubridate' was built under R version 4.1.3
```

```r
library(knitr)
```

```r
cdc <- read.csv('CDC-spotify.csv')
```

```r
cdc = cdc |>
  rename('Available Markets' = Available.Markets, 'Duration (sec)' = Duration..sec., 'Track Name' = Tra
```

```r
spotify_data = cdc |>
  select(-X)
```
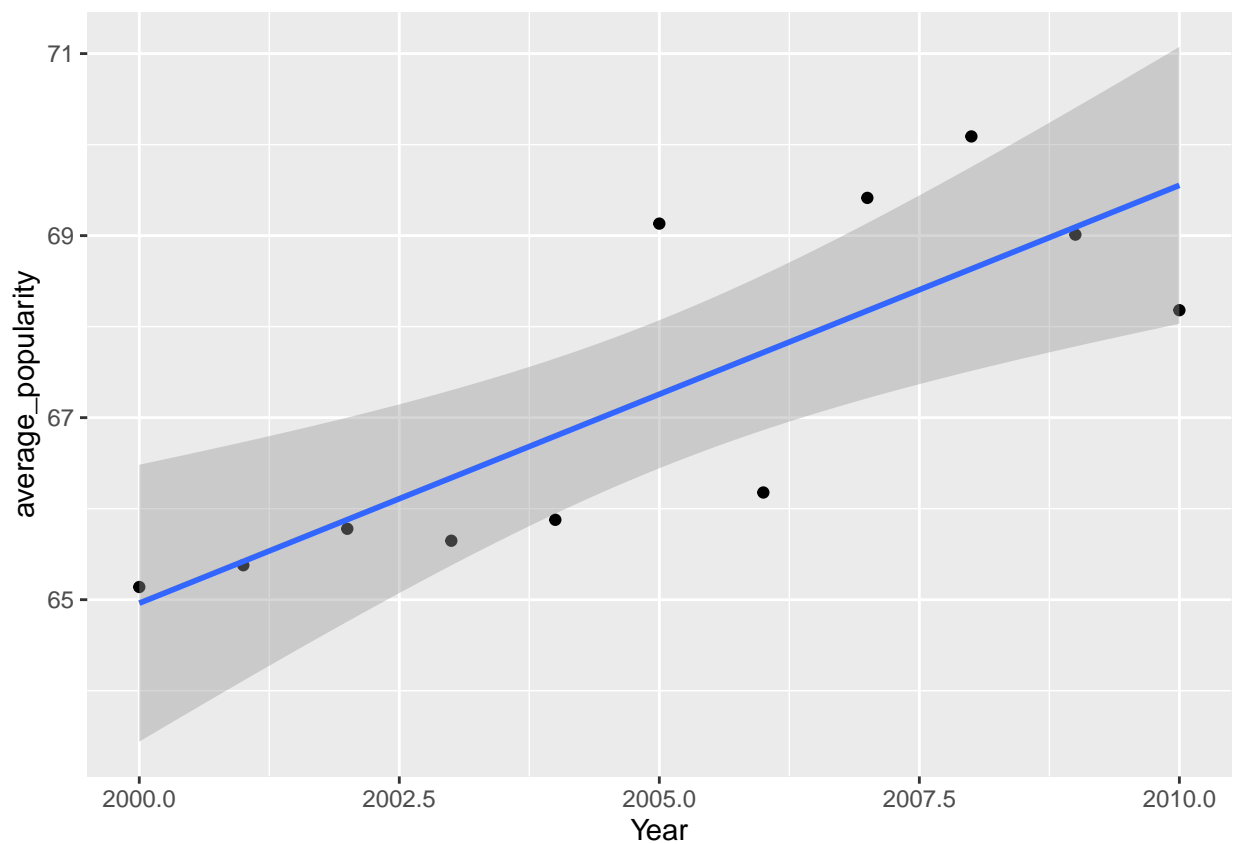
```r
average_years = spotify_data |>
  group_by(Year) |>
  summarize(average_popularity = mean(Popularity))
```

```r
average_years
```

```
## # A tibble: 11 x 2
##     Year average_popularity
##    <int>           <dbl>
##  1  2000            65.1
##  2  2001            65.4
##  3  2002            65.8
##  4  2003            65.6
##  5  2004            65.9
##  6  2005            69.1
##  7  2006            66.2
##  8  2007            69.4
##  9  2008            70.1
## 10  2009            69.0
## 11  2010            68.2
```
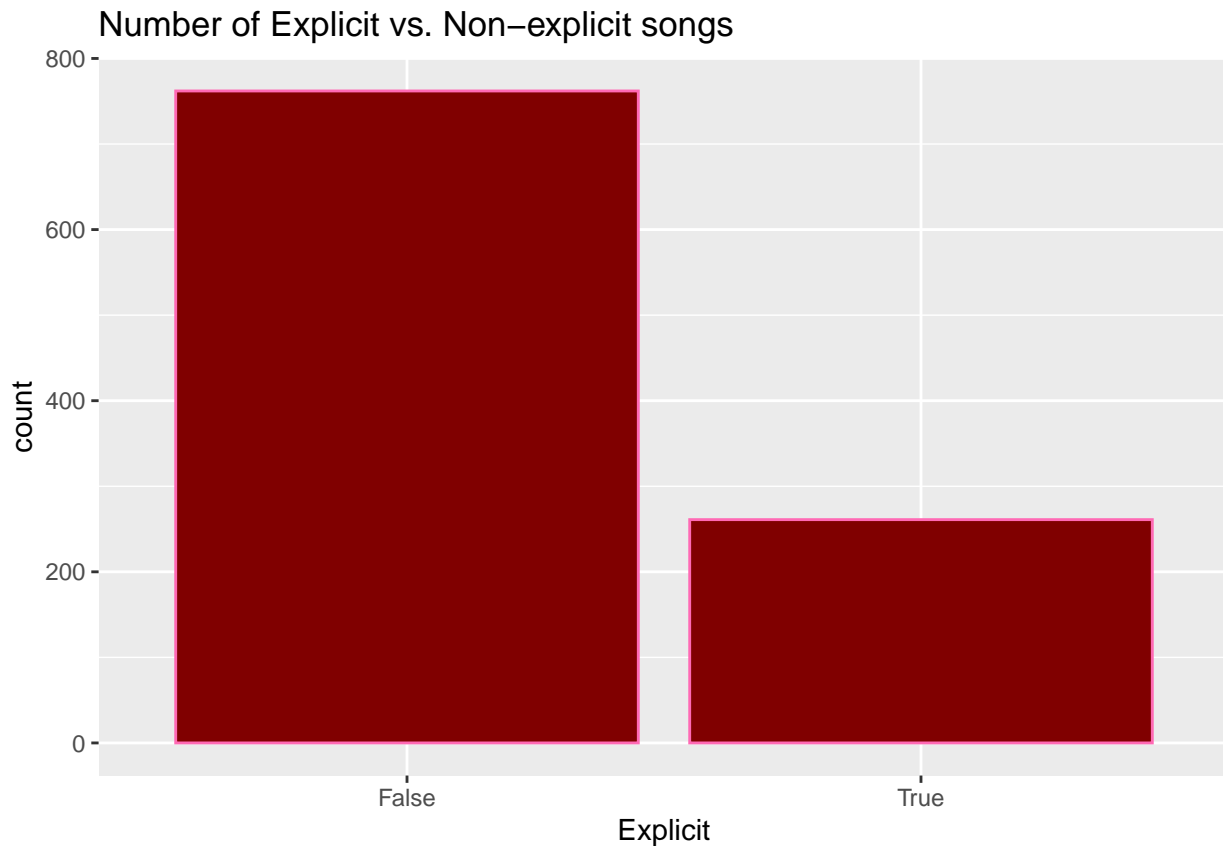
```
average_years |>
  ggplot(mapping = aes(x= Year, y = average_popularity)) +
  geom_point() +
  geom_smooth(method = "lm", na.rm = TRUE)
```

```
## 'geom_smooth()' using formula = 'y ~ x'
```



```
spotify_data |>
  ggplot(mapping = aes(x = Explicit)) +
```

```
geom_bar(color = "#FF69B4", fill = "#800000") +
labs(title = "Number of Explicit vs. Non-explicit songs")
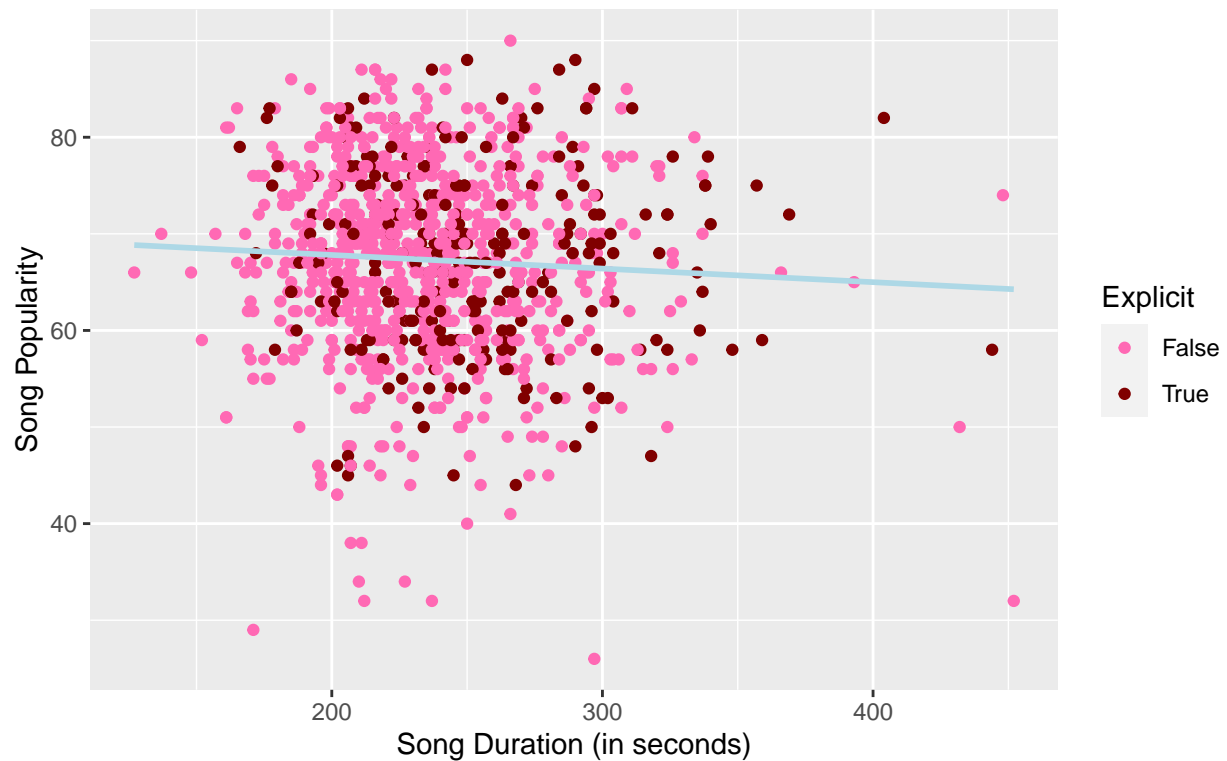```

## Number of Explicit vs. Non−explicit songs



```
average_explicit = spotify_data |>
  group_by(Explicit) |>
  summarize(Explicit_Popularity = mean(Popularity))
```

```
spotify_data |>
  ggplot(mapping = aes(x = `Duration (sec)`, y = Popularity)) +
  geom_point(aes(color = Explicit)) +
   scale_color_manual(values = c("#FF69B4", "#800000")) +
  geom_smooth(method = 'lm', color = "#ADD8E6", se = FALSE) +
    labs(title = "Relationship Between Song Popularity and Duration",
       x="Song Duration (in seconds)",
     y="Song Popularity",
       caption = "Source: Spotify API ")
```

```
## `geom_smooth()` using formula = 'y ~ x'
```

## Relationship Between Song Popularity and Duration



Source: Spotify API

```r
average_artists = spotify_data |>
  group_by(Artist) |>
  summarize(`Mean Popularity` = mean(Popularity),
            Count = n(),
            `Average Duration` = mean(`Duration (sec)`)) |>
  arrange(desc(Count))
```

```r
average_artists
```

```
## # A tibble: 475 x 4
##    Artist          `Mean Popularity` Count `Average Duration`
##    <chr>                       <dbl> <int>              <dbl>
##  1 Eminem                       74.5    18               290.
##  2 Rihanna                      75.4    15               237.
##  3 Kanye West                   76.9    14               235.
##  4 Britney Spears               71      13               212.
##  5 Beyoncé                      69.1    12               231.
##  6 USHER                        73.4    12               241.
##  7 Black Eyed Peas              67.8    11               261.
##  8 Alicia Keys                  65.6    10               268
##  9 50 Cent                      69.9     9               227.
## 10 Coldplay                     80.8     9               273.
## # i 465 more rows
```