

How your song becomes a hit

Group 4 - Caterina Bonomi, Nataliya Kharitonova, Ines Lindtner, Lucia Raffaelli

2023-02-07

Let's start with the basics!

With the advent of new technologies – especially the emergence of digital platforms – the music industry has been evolving and changing some of its structural features. Artists can now reach incredible success through the use of algorithmically made playlists on Spotify instead of having to promote an entire album through record labels. Clearly the dynamics have changed, and they keep changing, which makes this industry a very interesting landscape to analyze.

However, something that hasn't undergone particular changes is the use of charts to determine the most popular songs in different countries. Are charts really telling us which are the hit songs and what makes them so special? How does a song enter the top 10 in a chart? Why do songs have different survival rates?

This article aims at investigating which songs made it to the top of the charts over the years and what was their survival rate over time. We, as a team of marketing students – with a passion for music – decided to test this assumption with data of the top 200 charts from 2015-2021 gathered in the following 8 countries:

```
## [1] "au" "gb" "de" "ca" "us" "it" "fr" "jp"
```

For simplicity reasons, we decided to focus on data chosen from 3 out of the 8 countries - Australia, Italy, and the US – in the attempt of making a comprehensive analysis considering that these states are from completely different parts of the world.

All the songs present certain features, which can be divided into sound features (9) and technical features (4). Analyzing these variables allows us to get a feeling for the song and compare the scores of hit songs to other less successful songs. Here are the available features listed:

```
##           type           feature
## 1      Sound features    acousticness
## 2      Sound features    danceability
## 3      Sound features         energy
## 4      Sound features instrumentalness
## 5      Sound features         liveness
## 6      Sound features         loudness
## 7      Sound features    speechiness
## 8      Sound features         tempo
## 9      Sound features         valence
## 10 Technical features    time_signature
## 11 Technical features         mode
## 12 Technical features         key
## 13 Technical features    explicit
```

The first section of the article will present the different countries analyzed individually, whilst the second part will show a comparison among them.

The charts in the countries

Australia

We started our analysis by looking at the six most popular songs in Australia between 2015-2021, with their respective artist names, streams, and weighted ranking position. We calculated a weighted ranking (wgt_rank) in order to be able to compare songs which have stayed in the top charts for short periods of time to songs that have stayed longer but in lower positions.

##	trackName	artistName
## 3904	Shape of You	Ed Sheeran
## 979	One Dance	Drake feat. WizKid feat. Kyla
## 3647	Closer	The Chainsmokers feat. Halsey
## 2584	Starboy	The Weeknd feat. Daft Punk
## 2470	Lean On (feat. MØ & DJ Snake)	Major Lazer feat. MØ feat. DJ Snake
## 1483	God's Plan	Drake

##	wgt_rank
## 3904	139.64680
## 979	118.61555
## 3647	113.90544
## 2584	105.14068
## 2470	98.19961
## 1483	89.62535

After showing the most popular songs in Australia during this time frame, we decided to investigate what made these songs so successful. For this purpose, we chose to run a logarithmic regression, which only takes into consideration songs that were released after 2015. We made this choice because we were provided with data on charts starting from 2015, therefore considering songs that were released before that year would mean not analyzing the entire lifetime of the songs, hence making it impossible to fully estimate the true popularity that these songs had.

```
##
## Call:
## lm(formula = log(wgt_rank * 100 + 1) ~ log(danceability * 100 +
##      1) + log(speechiness * 100 + 1) + log(acousticness * 100 +
##      1) + log(energy * 100 + 1) + log(instrumentalness * 100 +
##      1) + log(liveness * 100 + 1) + log(loudness * 100 + 1) +
##      log(tempo * 100 + 1) + log(valence * 100 + 1) + time_signature +
##      mode + key + explicit + duration, data = chart_au_year)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.6143 -1.7691 -0.2952  1.4760  6.1996
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    -1.638e+01  3.503e+00  -4.675 3.08e-06 ***
## log(danceability * 100 + 1)  1.205e+00  1.738e-01   6.928 5.26e-12 ***
## log(speechiness * 100 + 1)   8.287e-03  6.647e-02   0.125  0.90079
## log(acousticness * 100 + 1)  4.786e-04  3.614e-02   0.013  0.98943
## log(energy * 100 + 1)      -5.377e-01  2.087e-01  -2.576  0.01003 *
## log(instrumentalness * 100 + 1) -2.454e-01  5.058e-02  -4.851 1.29e-06 ***
## log(liveness * 100 + 1)     -2.033e-01  7.030e-02  -2.892  0.00386 **
## log(loudness * 100 + 1)      4.099e+00  8.619e-01   4.755 2.08e-06 ***
## log(tempo * 100 + 1)        -2.350e-01  1.726e-01  -1.361  0.17348
## log(valence * 100 + 1)      -5.449e-02  7.986e-02  -0.682  0.49508
## time_signature      1.530e-01  1.589e-01   0.963  0.33566
## mode                -1.326e-01  8.272e-02  -1.603  0.10896
## key                  2.628e-03  1.110e-02   0.237  0.81285
## explicit             4.507e-02  9.654e-02   0.467  0.64066
## duration            -1.027e-06  8.904e-07  -1.154  0.24871
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.11 on 2793 degrees of freedom
## Multiple R-squared:  0.05632,    Adjusted R-squared:  0.05159
## F-statistic: 11.91 on 14 and 2793 DF,  p-value: < 2.2e-16
```

We begin the interpretation of the statistical output by assessing the model fit and we can first of all see that the model is significant, given the p-value lower than the usual significance level of 0.05. However, according to the multiple R-squared, the model explains only 5.6% of the variation of the dependent variable.

Looking at the coefficients, we can observe the following:

- Danceability, energy, instrumentalness, liveness and loudness are the only significant variables with a p-value lower than 0.05, meaning that we can favor the hypothesis that there is a non-zero correlation between these variables and the dependent variable weighted rank.
- Among the significant variables we can distinguish between the ones that positively affect the ranking position of the song – which are danceability and loudness – and those who negatively affect it – which are energy, instrumentalness and liveness.

Given these considerations, we can infer that Australians prefer songs that are somewhat more danceable and particularly loud (coefficients equal to 1.205 vs 4.099). Instead, songs that are somewhat more energetic (-0.5377), instrumental (-0.2454) and live(-0.2033) songs are less likely to become very popular.

```
##                                feols(log(wgt_ran..
## Dependent Var.:              log(wgt_rank*100+1)
##
## log(danceability x 100+1)    1.260*** (0.0614)
## log(speechiness x 100+1)     0.0218 (0.0654)
## log(acousticness x 100+1)    0.0014 (0.0307)
## log(energy x 100+1)         -0.5763* (0.1394)
## log(instrumentalness x 100+1) -0.2562* (0.0666)
## log(liveness x 100+1)       -0.2156 (0.0922)
## log(loudness x 100+1)       4.249** (0.6074)
## log(tempo x 100+1)         -0.2268 (0.1263)
## log(valence x 100+1)        -0.0655 (0.0842)
## time_signature              0.1515* (0.0450)
## mode                       -0.1265 (0.1162)
## key                        0.0023 (0.0122)
## explicit                   0.0575 (0.0960)
## duration                   -1.22e-6 (9.44e-7)
## Fixed-Effects:             -----
## year                       Yes
## _____
## S.E.: Clustered              by: year
## Observations                2,808
## R2                          0.05899
## Within R2                   0.05847
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Following this initial analysis, we decided to run another logarithmic regression using the release year as fixed effect. Indeed, the fixed effects logistic regression models have the ability to control for all songs' fixed characteristics (time independent). This allows us to understand how release year affects the ranking position eliminating other potential confounding effects on our dependent variable.

Looking at the output, we can see that different variables significantly affect the ranking position. Other than danceability, energy, instrumentalness and loudness, this regression shows us that time signature is also affecting the relationship. This is an interesting insight as it proves our assumption that the release year does compromise the analysis and therefore it does make sense to control for time invariant characteristics of the songs.

```
##                                feols(log(wgt_ran..
## Dependent Var.:              log(wgt_rank*100+1)
##
## log(danceability x 100+1)      0.6658* (0.1804)
## log(speechiness x 100+1)       0.0525 (0.0448)
## log(acousticness x 100+1)     -0.0288 (0.0870)
## log(energy x 100+1)           0.0212 (0.1119)
## log(instrumentalness x 100+1) -0.0394 (0.1856)
## log(liveness x 100+1)         -0.0610 (0.1218)
## log(loudness x 100+1)         2.816* (0.7674)
## log(tempo x 100+1)            -0.1389 (0.0800)
## log(valence x 100+1)          0.1301 (0.1887)
## time_signature                 0.1750 (0.1401)
## mode                          -0.1708 (0.1578)
## key                           -0.0058 (0.0139)
## explicit                      -0.1469 (0.1699)
## duration                      7.51e-7 (8.67e-7)
## Fixed-Effects:                -----
## year                          Yes
## artistName                    Yes
## _____
## S.E.: Clustered                by: year
## Observations                  2,808
## R2                            0.63523
## Within R2                     0.03178
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Moreover, we also decided to run a third logarithmic regression looking at how the release year and the role of the artist separately affect the ranking position. With respect to the previous analysis, in this case only danceability and loudness appear to be significantly and positively affecting the ranking position. We can therefore understand that the role of the artist affects the relationship between the ranking position and the independent variables. We can interpret this in two ways:

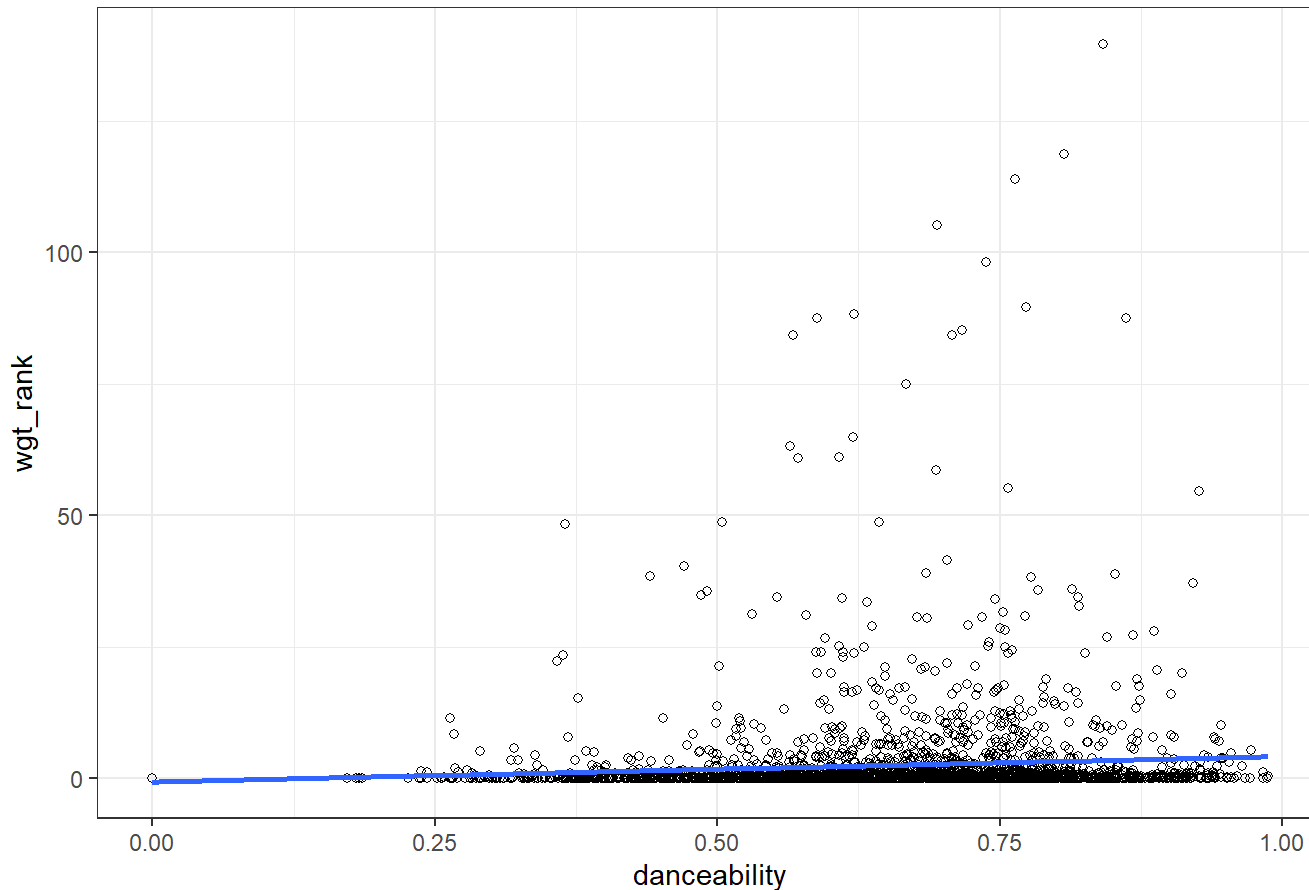
- The artist role affects the features of the song, which in turn have a low effect on the overall score.
- The artist role affects the score, which indicates that his/her stardom has the biggest influence on the score.

The conclusion that we can draw from these assumptions is that the more famous the person is, the more likely the song is to become a hit. This model, when accounting for the release year and the artist explains 63.5% of the data.

At this point of the analysis, we chose to create a series of graphical representations to better show the existing correlations between some of the above-mentioned significant variables and the ranking position. In particular we chose danceability, loudness and energy because they resulted to be the most significant in the previous regressions and we thought it would be interesting to compare them and their effects on the dependent variable.

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

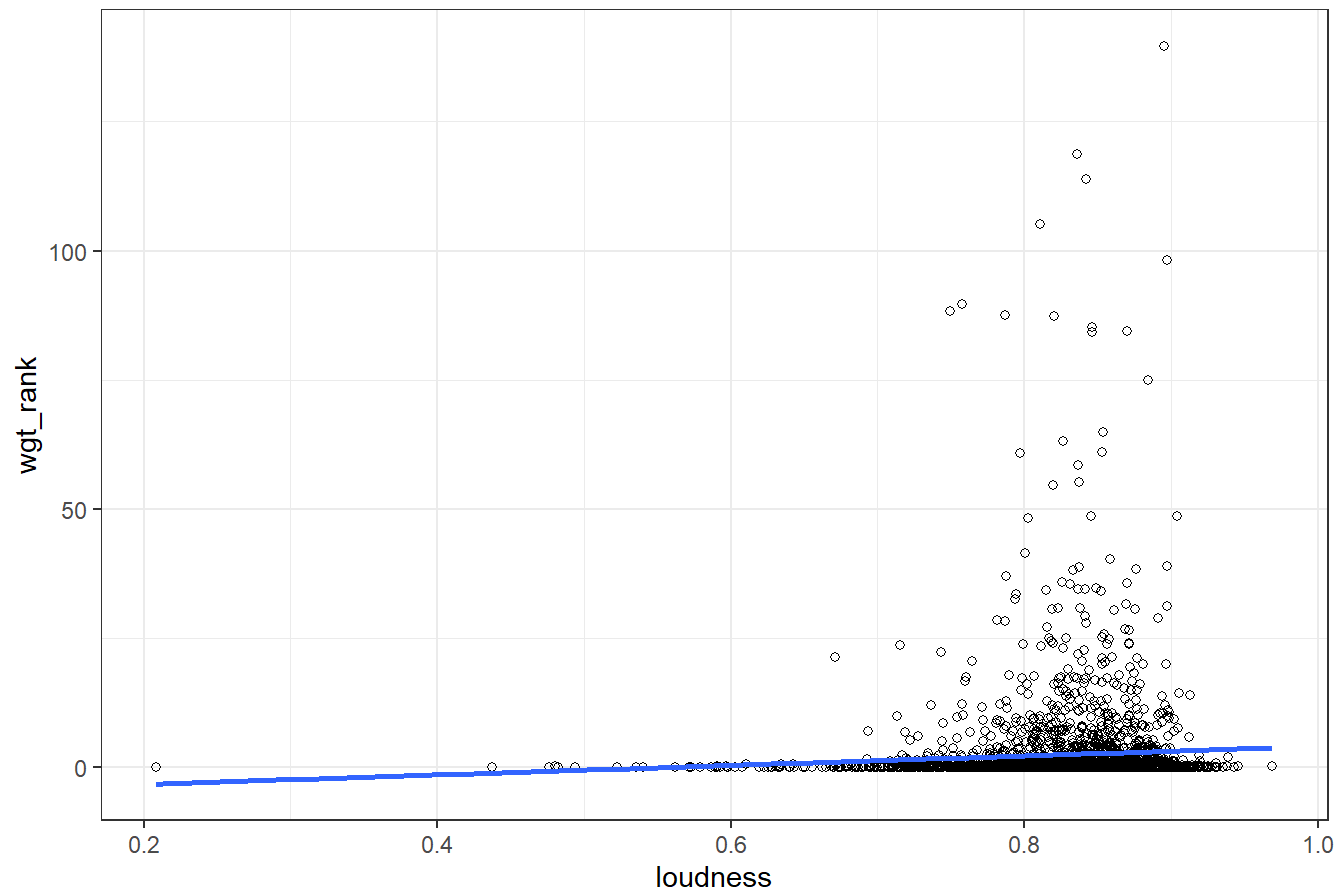
Linear Regression: Danceability & Ranking



Looking at the graph above, we can observe a linear correlation between rank position and danceability as the more danceable a song is, the more likely it is to obtain a higher ranking in the charts. However, the correlation is not that strong, which is represented by the almost flat slope of the line.

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

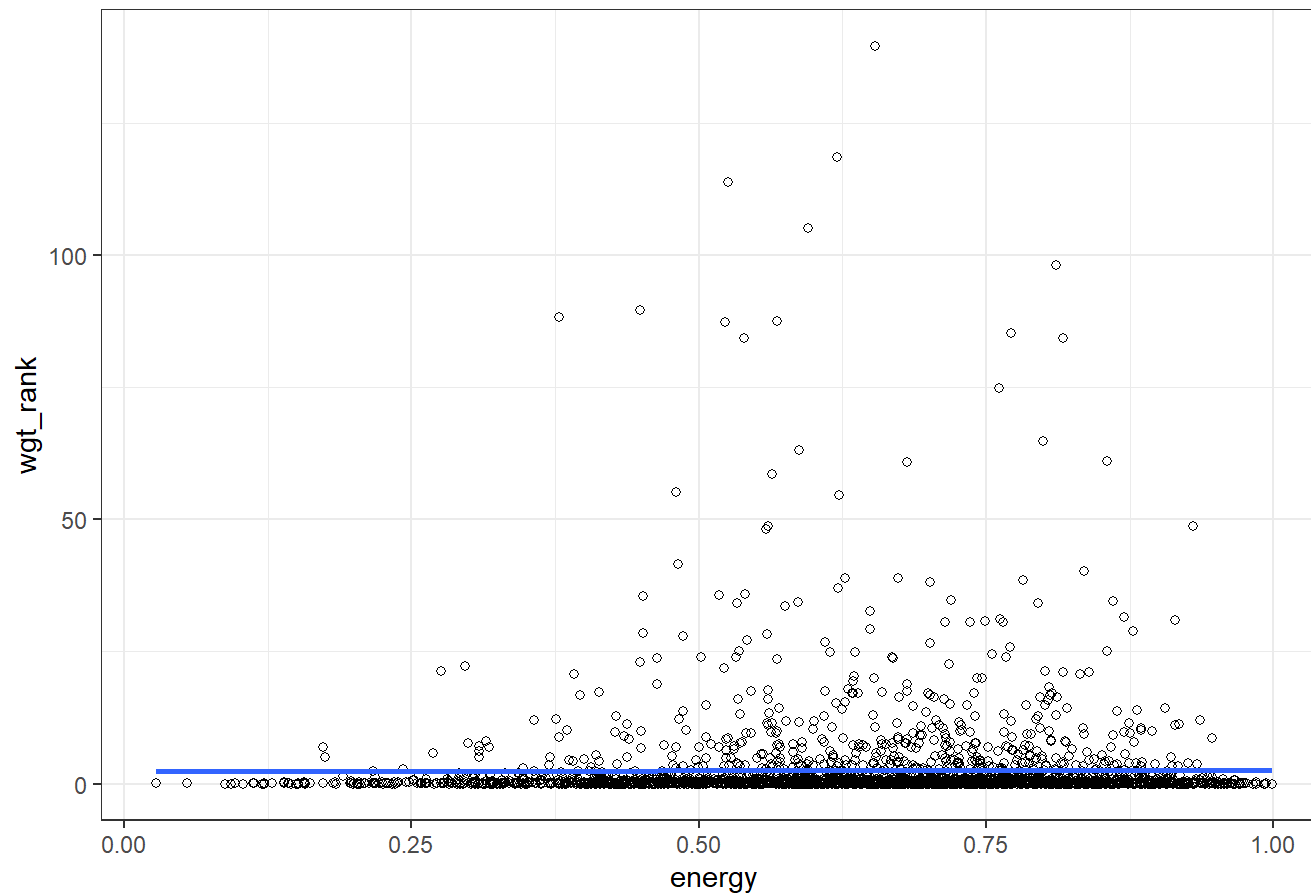

Linear Regression: Loudness & Ranking



Looking at the graphical representation of the relationship between ranking position and loudness, we can recognize again a linear correlation, this time stronger than the previous one. This goes in line with the regression results, as loudness had a higher coefficient indicating a stronger effect on the ranking position.

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

Linear Regression: Energy & Ranking



The graph above shows that there is no strong correlation between ranking position and energy. Indeed, the line connecting the observations is almost horizontal, clearly pointing out the absence of relationship between energy and the dependent variable.

After showing the existing or non-existing correlations between these three variables – danceability, loudness and energy – and the ranking position, we run a t-test investigating whether these variables taken individually have an effect on the top 10 selection of songs in Australia.

```
##
## Welch Two Sample t-test
##
## data: log(danceability * 100 + 1) by TopTen
## t = -3.9424, df = 9.2919, p-value = 0.003188
## alternative hypothesis: true difference in means between group 0 and group 1 is not equal to 0
## 95 percent confidence interval:
## -0.24099467 -0.06579754
## sample estimates:
## mean in group 0 mean in group 1
## 4.158482 4.311878
```

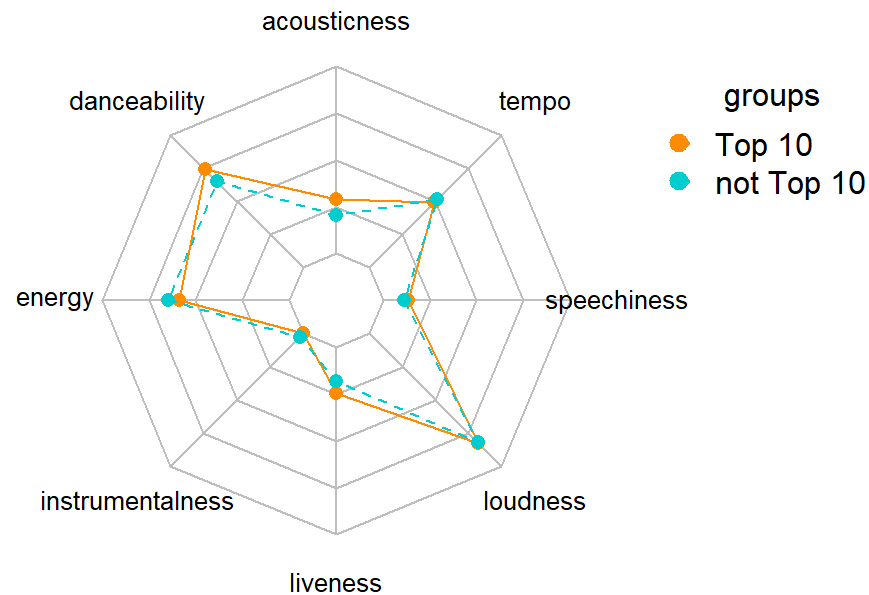
The first t-test aims at analyzing the effect of danceability on the population of interest – top 10 songs in Australia. Given that the p-value is lower than 0.05, we can reject the null hypothesis and accept the alternative one. This means that the difference in means is not equal to 0, showing that it is very likely that danceability affects the top 10 songs in Australian charts. The conclusion we can get from this is that generally the best 10 songs are more danceable.

```
##
## Welch Two Sample t-test
##
## data: log(loudness * 100 + 1) by TopTen
## t = -0.11338, df = 9.1017, p-value = 0.9122
## alternative hypothesis: true difference in means between group 0 and group 1 is not equal to 0
## 95 percent confidence interval:
## -0.04561335 0.04125229
## sample estimates:
## mean in group 0 mean in group 1
## 4.420301 4.422482
```

```
##  
## Welch Two Sample t-test  
##  
## data: log(energy * 100 + 1) by TopTen  
## t = 0.97016, df = 9.1441, p-value = 0.3569  
## alternative hypothesis: true difference in means between group 0 and group 1 is not equal to 0  
## 95 percent confidence interval:  
## -0.09270671 0.23252083  
## sample estimates:  
## mean in group 0 mean in group 1  
## 4.140599 4.070692
```

The other two t-test we run resulted to be non-significant as the p-values are higher than 0.05 in both cases. The conclusion is that energy and loudness do not affect the ranking of the first 10 songs, which therefore do not necessarily need to be energetic and loud in order to be ranked highly.

Spiderplot of features in Australia

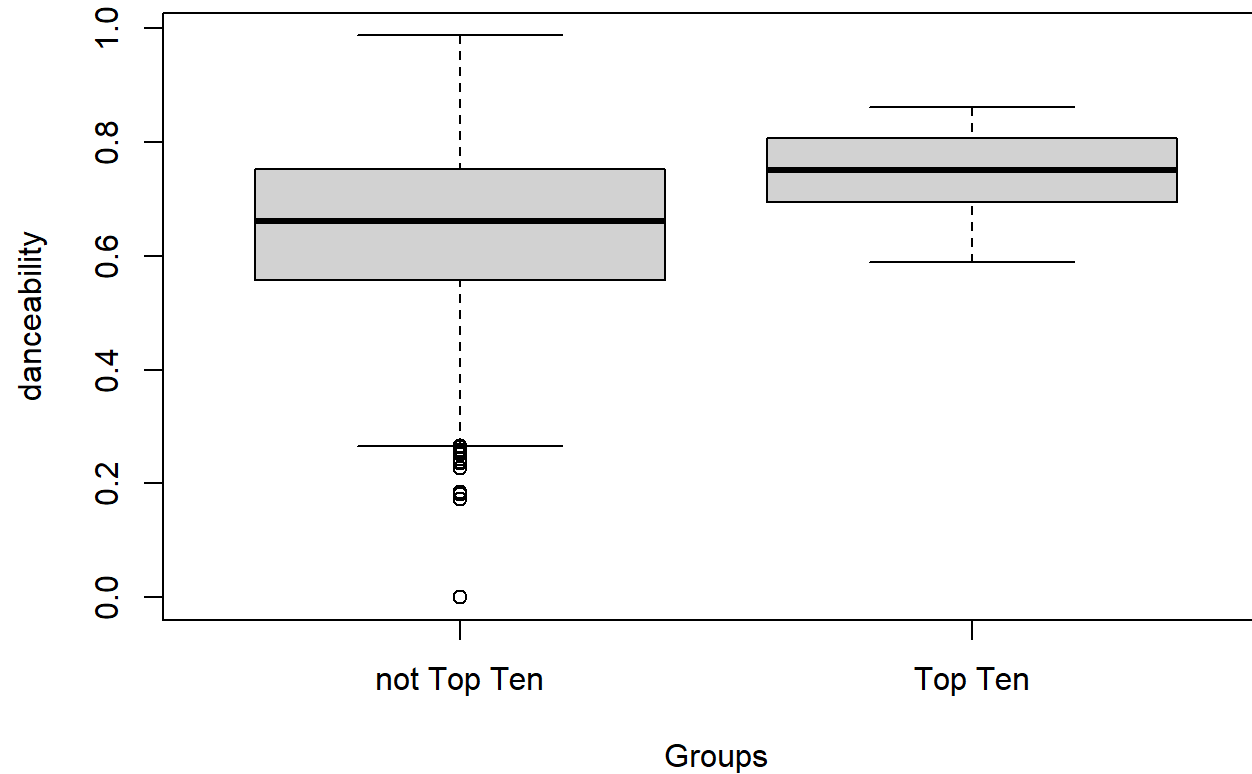


To get a better visual understanding of the features and their importance to the ranking position of the songs, we conducted the radar plot with 2 groups: Top 10 tracks in Australia and other songs. As we can see at a first glance, the results look similar for both clusters. Looking more closely, it is noticeable that not Top 10 songs are generally more energetic and faster in terms of pace, while top 10 songs are more danceable and acoustic tracks, which are then more likely to gain the first positions in the Australian charts. One interesting insight is that the live performed songs are also more likely to become “Top 10 tracks”. However, according to the linear regression results “liveness” leads to a lower weighted rating and, therefore, is connected to less likeable songs.

Moreover, both groups are quite low on instrumentalness, acousticness, speechiness, indicating that these features are not very typical features for songs in the charts. On the other hand, danceability, loudness and energy have high values overall.

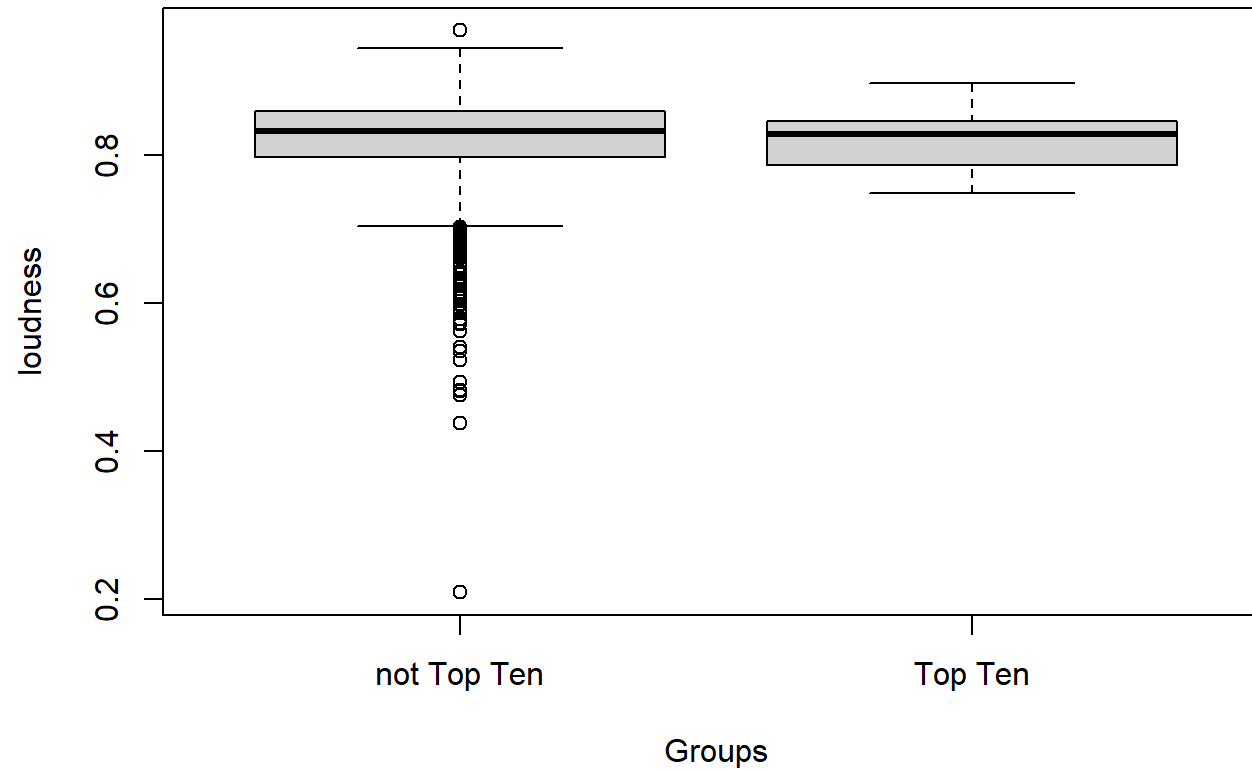
To further investigate the effect of the three variables showing the highest values – danceability, loudness and energy – we decided to draw some boxplots comparing top10 songs vs not top10 songs.

Boxplot - Danceability



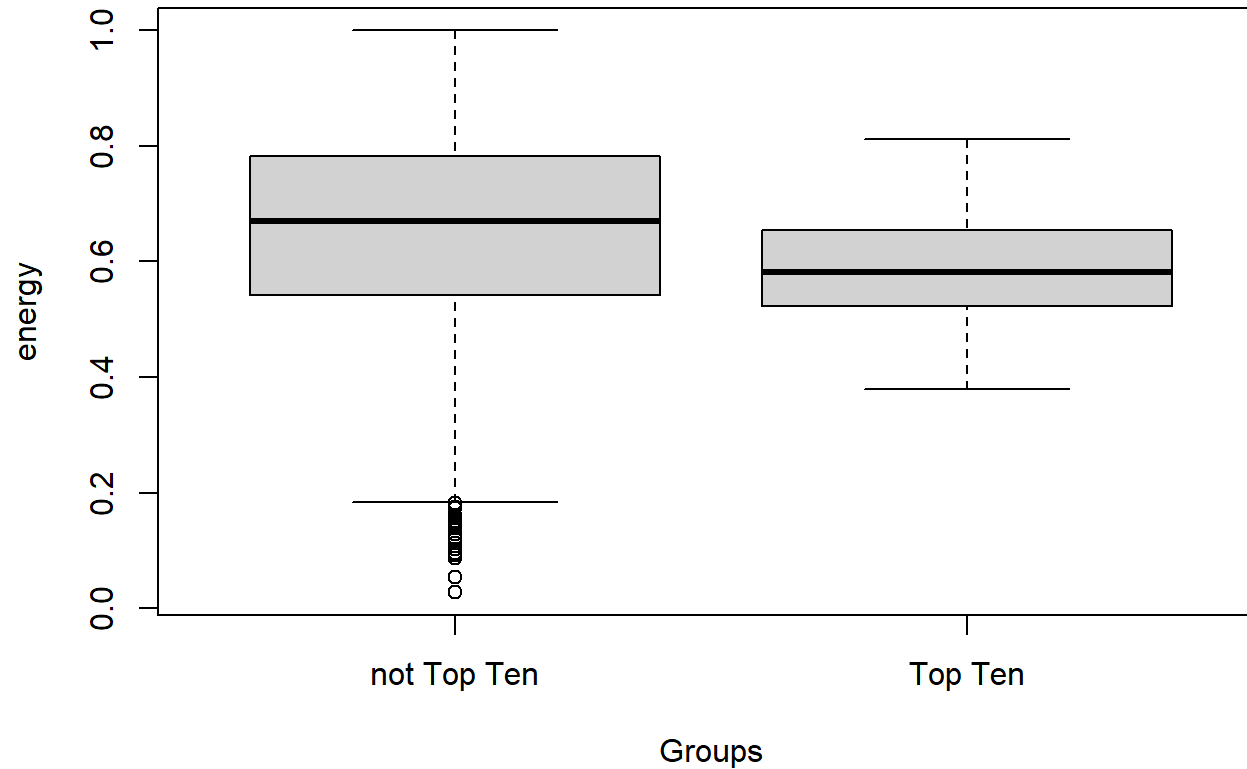
This boxplot shows the difference between top10 and non top10 songs in relation to danceability. We start by first comparing the medians, which are different with the median of the top10 songs being higher. This tells us that top10 songs are generally more danceable. Looking at the ranges, the non top10 songs present several outliers and the dispersion is bigger, which we need to take into consideration when drawing the final conclusion. Overall, we can say that danceability does indeed impact the ranking of songs and is generally more present in songs that reach the top10 charts.

Boxplot - Loudness



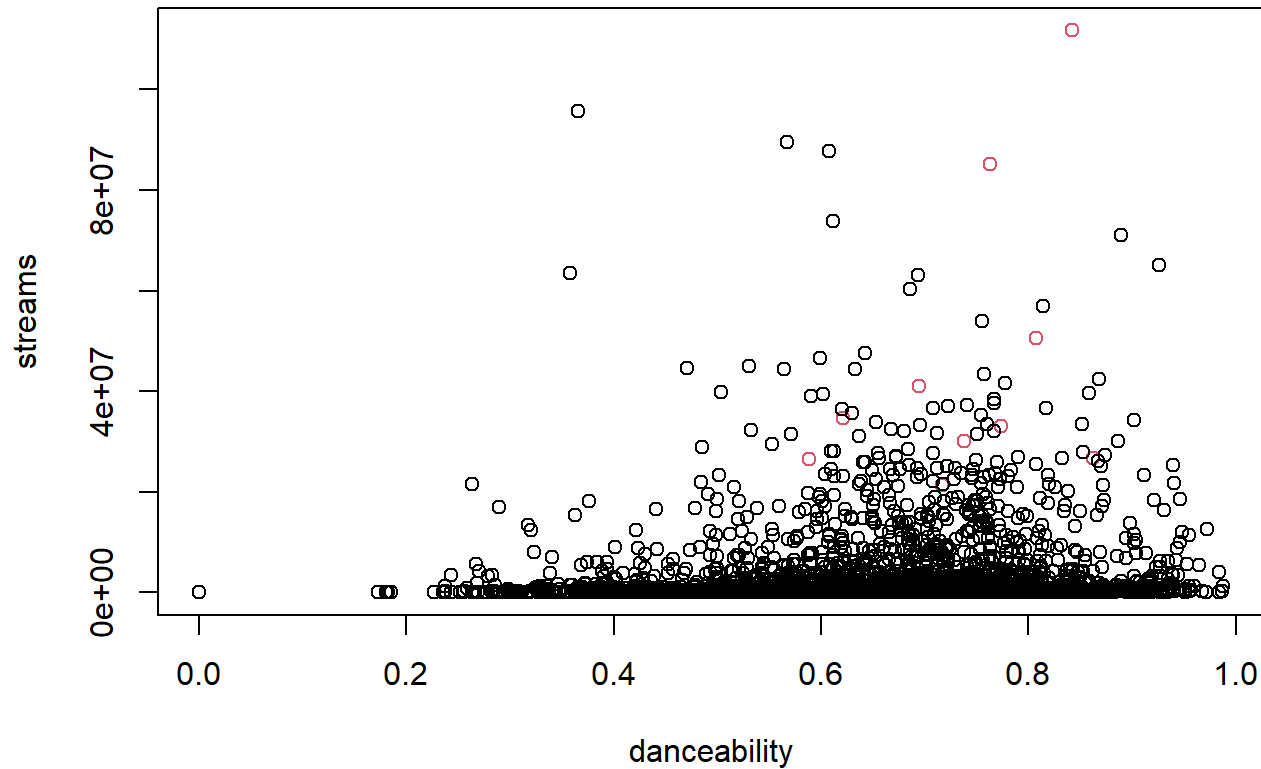
The second boxplot shows the difference between top10 and non top10 songs in relation to loudness. In this case the medians do not differ much, indicating that loudness does not make such difference between the two populations of songs. We do however observe many outliers in the not top10 group, which could distort our final insights.

Boxplot - Energy



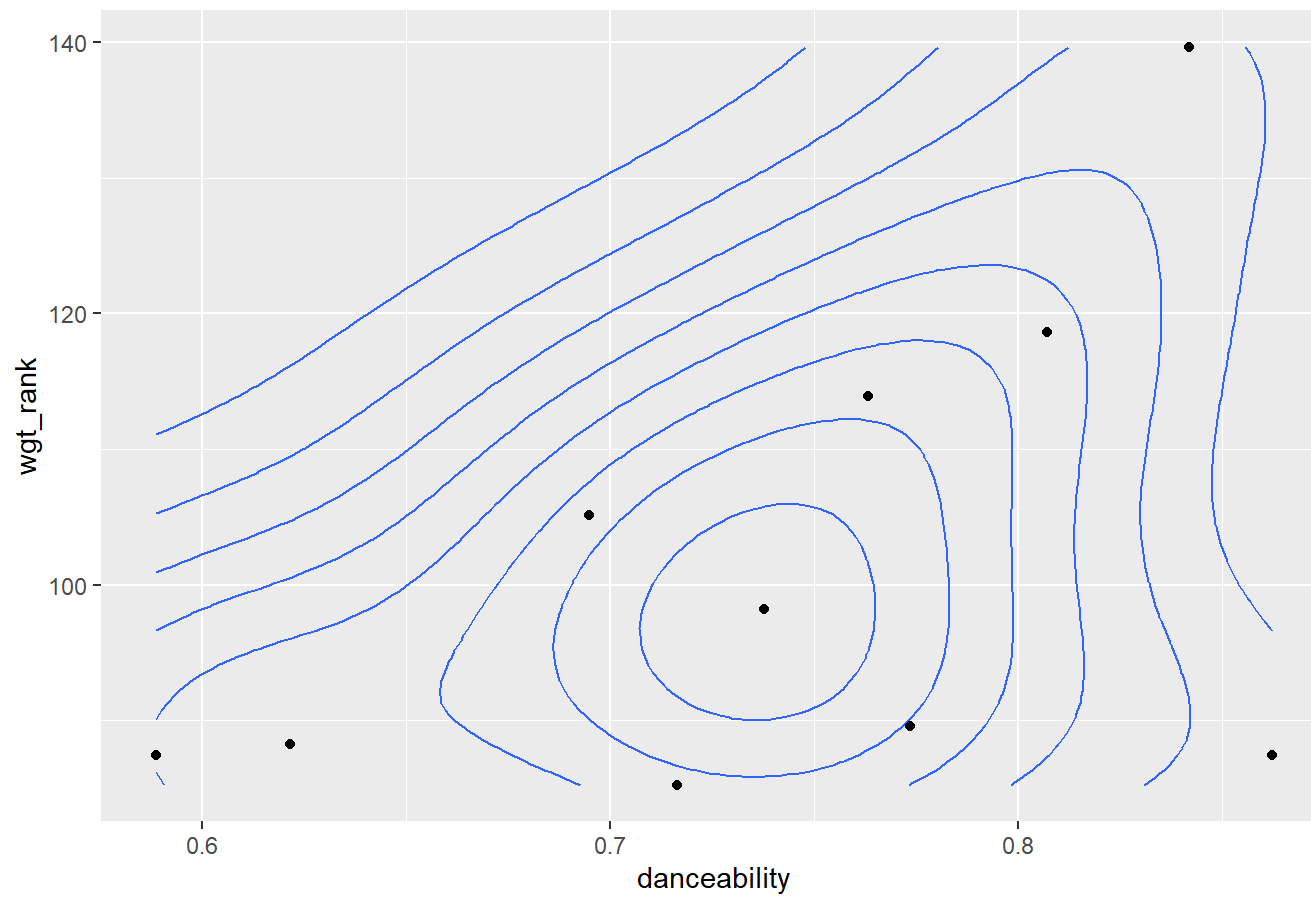
Regarding the last boxplot, we can see that the medians differ with the one of non top10 songs being slightly higher than top10. This indicates that energy is actually an indicator of less likeable songs, even though in this case we also see few outliers so we need to consider that results may be distorted. When considering the t-test conducted before we find that although there is a visible difference, this difference is statistically insignificant.

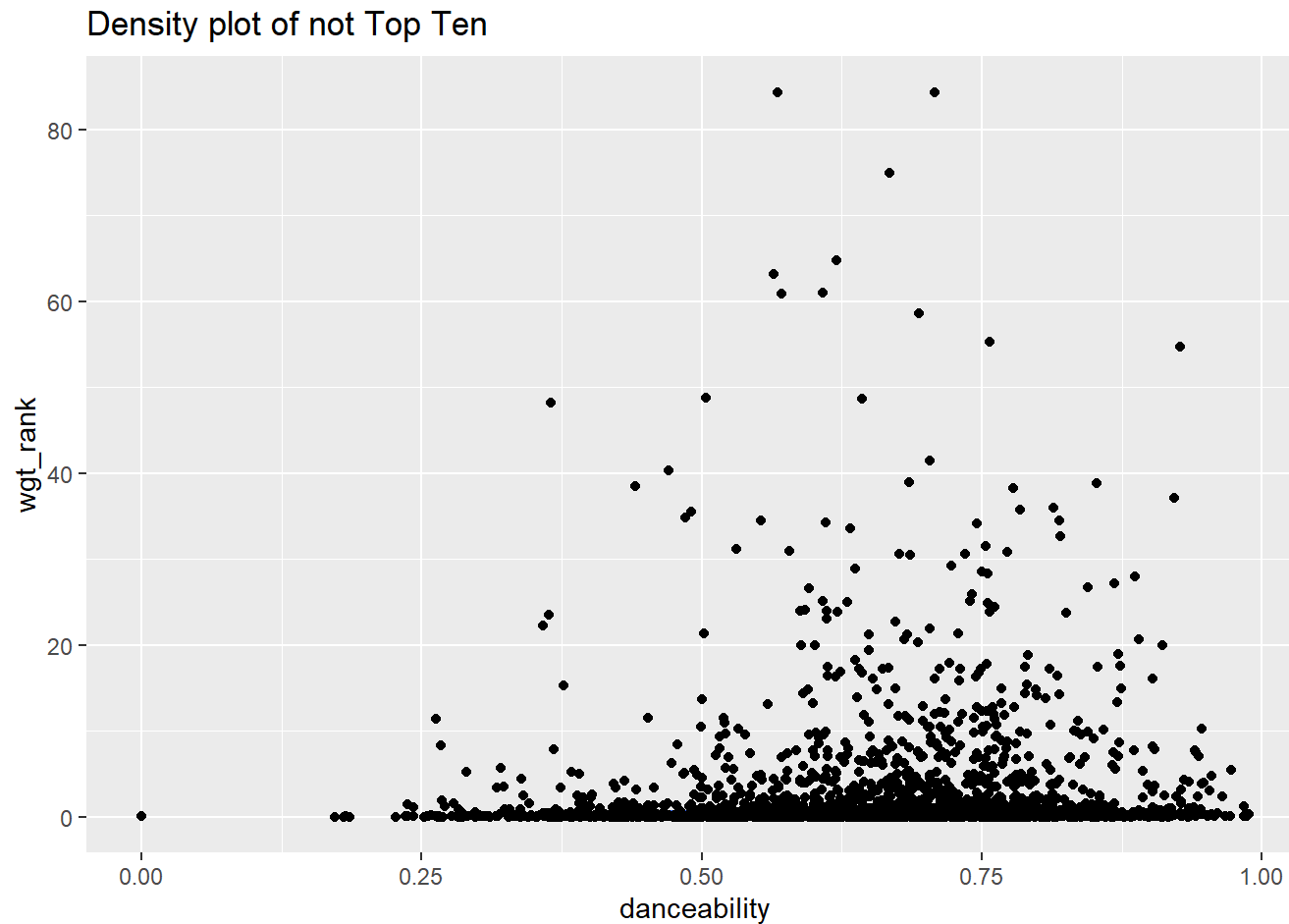
Scatterplot: Danceability & Streams



We then decided to graph the relationship between streams and danceability in a scatterplot, as the latter resulted to be the most significant variable. From this graph we can observe that when danceability is between 0.5 and 0.8 the number of streams increases. Therefore, we can assume that there is also a correlation between danceability and streams. In this graph the top ten songs are highlighted by being marked red.

Density plot of the Top Ten

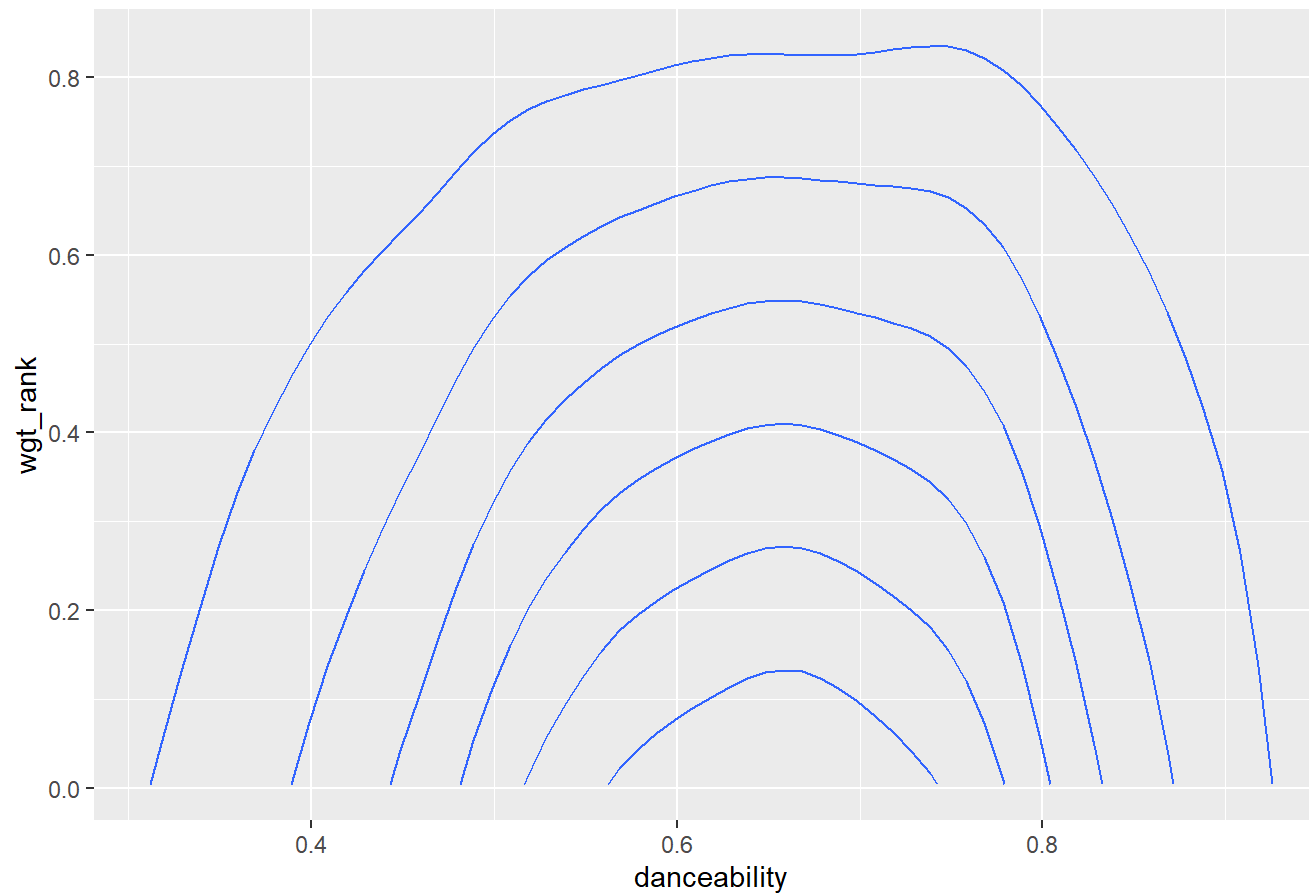




Here again we look at the relationship between danceability and weighted ranking through density plots. The first one is presenting the data of the top ten songs, while the second one shows all other songs. We can see via the density lines that for danceability values between 0.6 and 0.8 the ranking position increases for the top ten, showing us once again that danceability is a very valued feature when it comes to higher ranked songs. The second graph cannot produce visible density lines due to the dense amount of data points close to zero. Because of this the density lines are not visible in this graph. It shows though, that a lot of highly ranked songs have a danceability between 0.5 and 0.75.

If we want to take a look at the density lines we can plot the data with a focus on the density:

Density plot of not Top Ten



This shows, what we already saw before, that there is a lot of weight at the bottom.

Italy

We will now examine the data from Italy. As we did previously for Australia, we present the most popular songs in Italy between 2015 and 2021, along with their respective artist names, streams, and weighted ranking position. The outcomes are as follows:

##	trackName	artistName
## 3976	Shape of You	Ed Sheeran
## 2206	Despacito (Featuring Daddy Yankee)	Luis Fonsi feat. Daddy Yankee
## 2435	Let Me Love You	DJ Snake feat. Justin Bieber
## 500	Senza Pagare VS T-Pain	J-AX feat. Fedez feat. T-Pain
## 2950	Roma - Bangkok	Baby K feat. Giusy Ferreri
## 2550	Lean On (feat. MØ & DJ Snake)	Major Lazer feat. MØ feat. DJ Snake

##	streams	wgt_rank
## 3976	50395978	108.63505
## 2206	31203635	90.72688
## 2435	19862125	90.52313
## 500	28044304	89.97121
## 2950	17065561	89.27414
## 2550	20733108	88.60017

The question that arises once more is what these songs have in common, or more specifically, what makes these songs popular in Italy. We decided to perform another logarithmic regression with songs released after 2015.

```
##
## Call:
## lm(formula = log(wgt_rank * 100 + 1) ~ log(danceability * 100 +
##      1) + log(speechiness * 100 + 1) + log(acousticness * 100 +
##      1) + log(energy * 100 + 1) + log(instrumentalness * 100 +
##      1) + log(liveness * 100 + 1) + log(loudness * 100 + 1) +
##      log(tempo * 100 + 1) + log(valence * 100 + 1) + time_signature +
##      mode + key + explicit + duration, data = chart_it_year)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.3878 -1.8245 -0.3362  1.5188  6.3925
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    -1.902e+01  3.221e+00  -5.904 3.89e-09 ***
## log(danceability * 100 + 1)  1.112e+00  1.667e-01   6.675 2.87e-11 ***
## log(speechiness * 100 + 1)   6.380e-02  5.672e-02   1.125 0.260743
## log(acousticness * 100 + 1)  2.957e-02  3.338e-02   0.886 0.375752
## log(energy * 100 + 1)      -2.051e-01  1.940e-01  -1.057 0.290498
## log(instrumentalness * 100 + 1) -1.723e-01  5.897e-02  -2.922 0.003497 **
## log(liveness * 100 + 1)     -1.752e-01  6.131e-02  -2.857 0.004300 **
## log(loudness * 100 + 1)      4.090e+00  8.184e-01   4.998 6.08e-07 ***
## log(tempo * 100 + 1)        -1.358e-01  1.611e-01  -0.843 0.399417
## log(valence * 100 + 1)       5.921e-03  7.355e-02   0.081 0.935837
## time_signature      2.918e-01  1.229e-01   2.376 0.017578 *
## mode                -6.217e-02  7.324e-02  -0.849 0.396054
## key                 -9.170e-05  9.859e-03  -0.009 0.992579
## explicit            -3.339e-01  8.831e-02  -3.780 0.000159 ***
## duration            -9.295e-07  8.249e-07  -1.127 0.259921
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.078 on 3431 degrees of freedom
## Multiple R-squared:  0.05194,    Adjusted R-squared:  0.04807
## F-statistic: 13.43 on 14 and 3431 DF,  p-value: < 2.2e-16
```

We begin analyzing the model fit. As evidenced by the low p-value (<0.05), the results are significant. However, the model only explains about 5.2% of the data variability, according to the multiple R-squared.

Looking at the coefficients, we can see the following:

- The only significant variables, with p-values lower than 0.05, are danceability, instrumentalness, liveness, loudness, time signature, and explicit.
- Among the significant variables, we can distinguish those that positively affect the song's ranking position - danceability, loudness, and time signature - and those that negatively affect it - instrumentalness, liveness, and explicit.

This allows us to state that Italians prefer songs that are high in danceability, loudness (respectively with coefficients of 1.112 and 4.090). Instead, songs that are high in instrumentalness, liveness (-0.1723 and -0.1752 respectively) and that are also explicit (-0.3339) are less likely to be hits in Italy.

```
##                                feols(log(wgt_ran..
## Dependent Var.:              log(wgt_rank*100+1)
##
## log(danceability x 100+1)      1.217** (0.1843)
## log(speechiness x 100+1)       0.0994 (0.0562)
## log(acousticness x 100+1)      0.0418 (0.0325)
## log(energy x 100+1)           -0.2161 (0.2082)
## log(instrumentalness x 100+1) -0.1818* (0.0550)
## log(liveness x 100+1)         -0.1935* (0.0414)
## log(loudness x 100+1)         4.071*** (0.1961)
## log(tempo x 100+1)            -0.1489 (0.1735)
## log(valence x 100+1)          -0.0090 (0.0964)
## time_signature                 0.2708* (0.0686)
## mode                          -0.0657 (0.0394)
## key                           -0.0007 (0.0089)
## explicit                      -0.3137 (0.1815)
## duration                      -1.26e-6. (5.04e-7)
## Fixed-Effects:                -----
## year                          Yes
## _____
## S.E.: Clustered                by: year
## Observations                   3,446
## R2                             0.05913
## Within R2                      0.05527
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Also in this case, we use a logarithmic regression with fixed effects to see if the release year influences the popularity of a song in a rank. The results show that, all the previously significant features besides explicit are significant when accounting for the release year. Furthermore, the within R-squared increased to 5.5%, hence a higher share of variability of the data can now be explained. This demonstrates that the track's release year compromised our previous results.

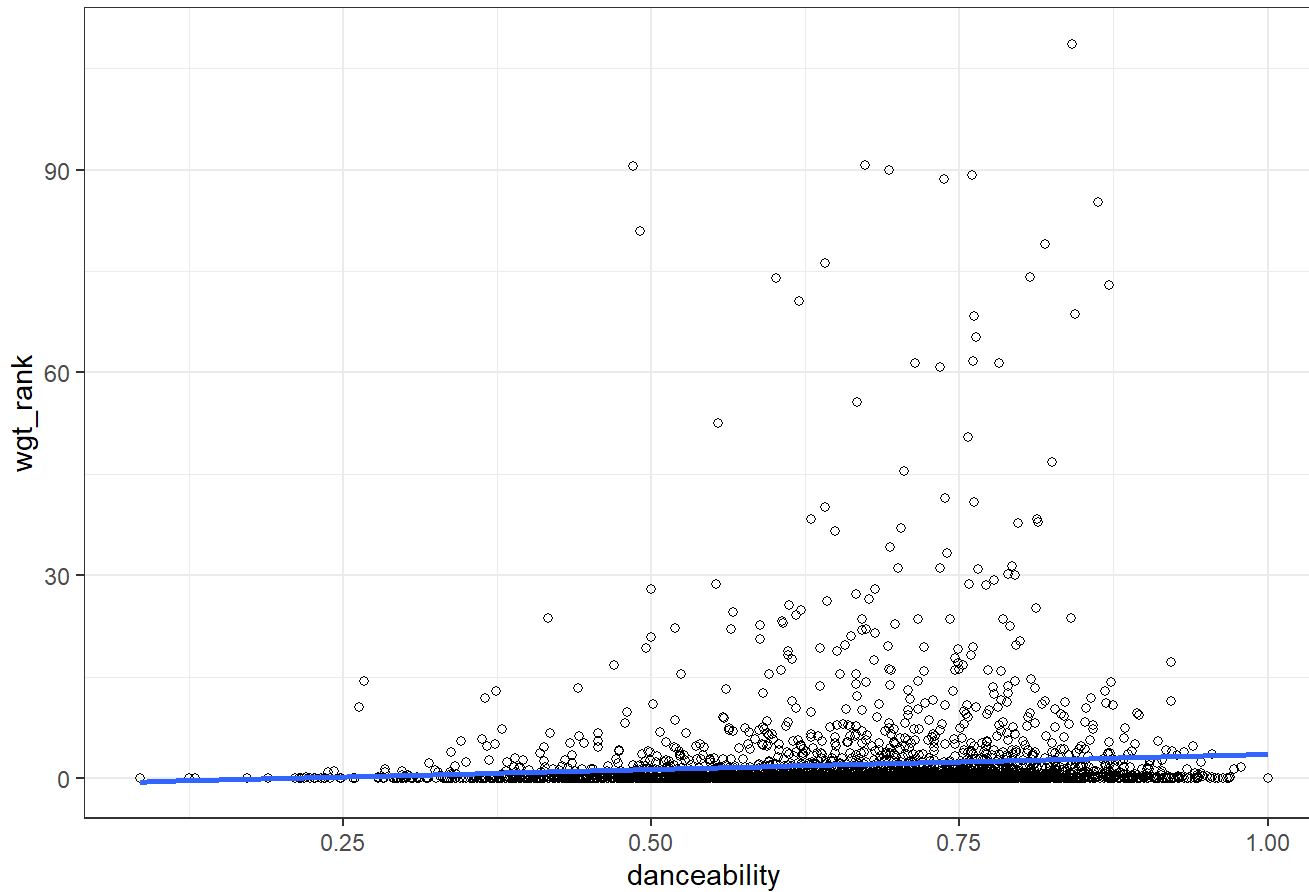
```
##                                feols(log(wgt_ran..
## Dependent Var.:              log(wgt_rank*100+1)
##
## log(danceability x 100+1)      0.6335* (0.1903)
## log(speechiness x 100+1)      -0.0124 (0.0414)
## log(acousticness x 100+1)     0.0145 (0.0194)
## log(energy x 100+1)           0.0351 (0.2288)
## log(instrumentalness x 100+1) -0.0464 (0.0687)
## log(liveness x 100+1)         -0.1972. (0.0790)
## log(loudness x 100+1)         3.216** (0.4342)
## log(tempo x 100+1)            -0.0826 (0.1094)
## log(valence x 100+1)          0.1222 (0.1185)
## time_signature                0.1573 (0.0729)
## mode                          0.0630 (0.0636)
## key                           -0.0014 (0.0160)
## explicit                      -0.0683 (0.1746)
## duration                      5.26e-7 (4.33e-7)
## Fixed-Effects:                -----
## year                          Yes
## artistName                    Yes
##
## _____
## S.E.: Clustered                by: year
## Observations                  3,446
## R2                           0.59535
## Within R2                     0.02853
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Again we run a logarithmic regression with two fixed effects: the release year and the artist's name. Also in this case the multiple R-squared increased tremendously and the model is now able to explain 59.5% of the variation in our data. In this regression the only significant results are danceability and loudness. We can infer that the artist's name affects the relationship between the song's rank and the independent variables.

We then proceed to do a graphical representation of the existing correlations between some of the previously mentioned variables and ranking position. We focused on danceability, loudness, and energy once more.


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

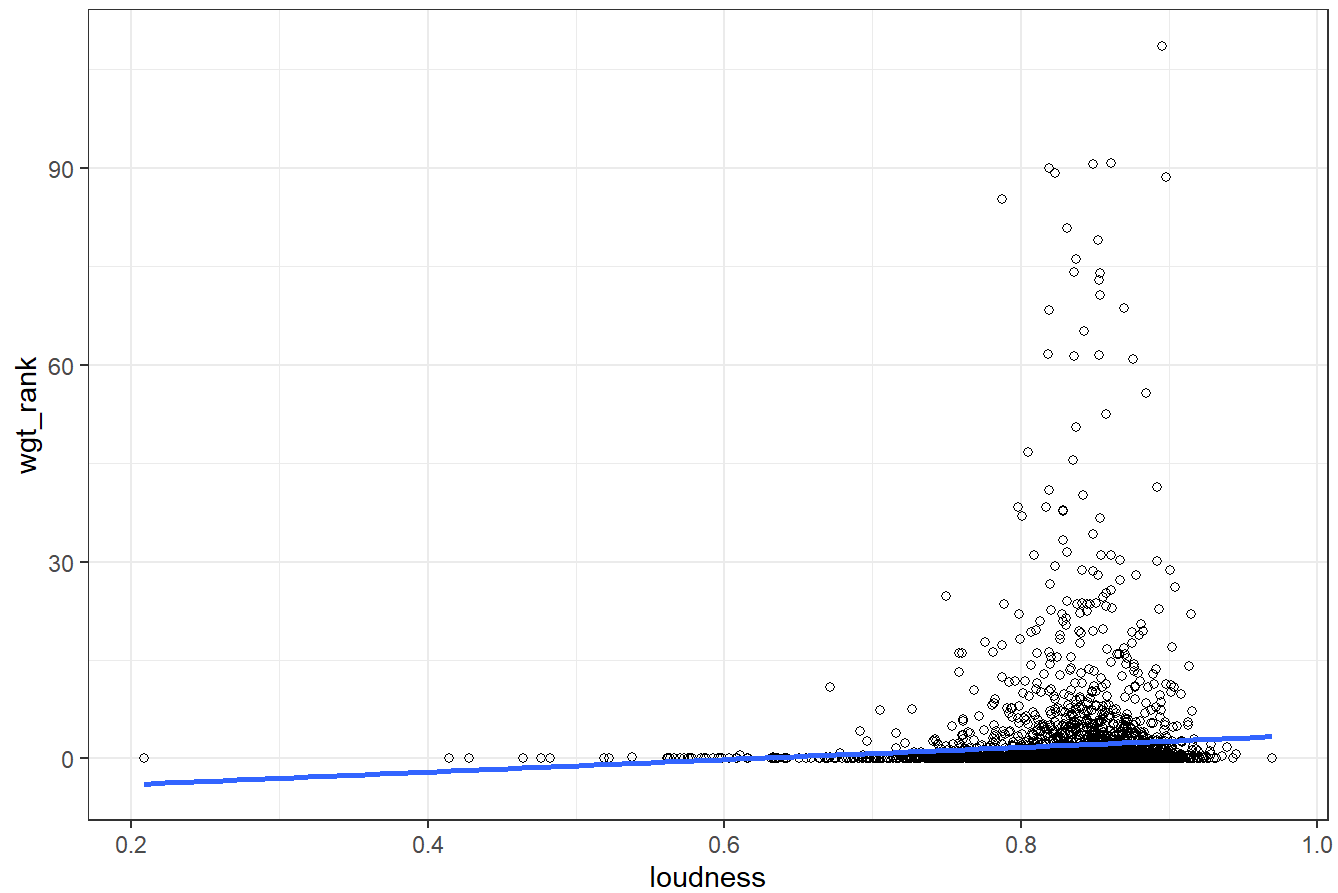
Linear Regression: Danceability & Ranking



We can identify a positive linear correlation between ranking position and danceability by looking at the graphical representation of the relationship between these two variables. The more suitable the track is for dancing, the better the chances of gaining a higher weighted ranking.

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

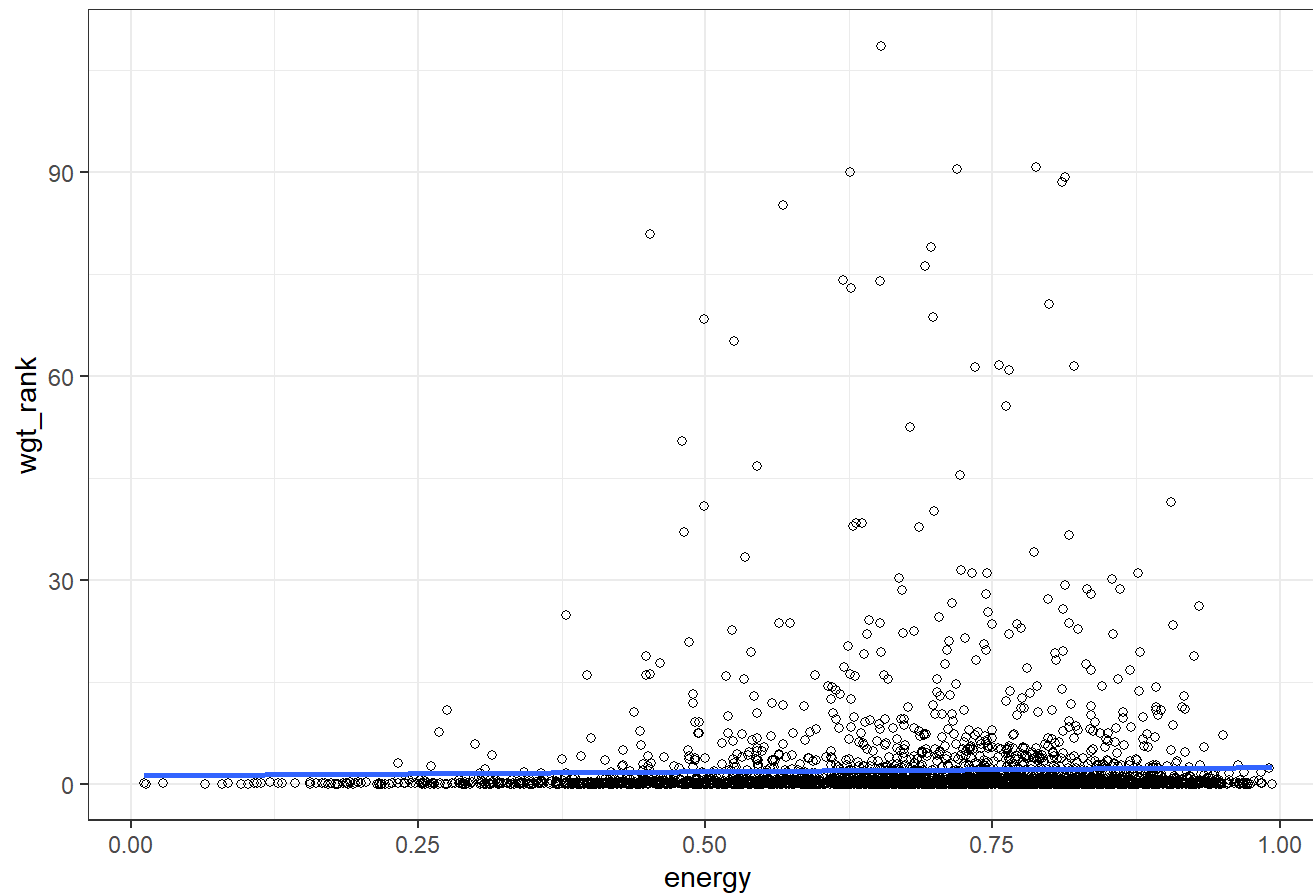
Linear Regression: Loudness & Ranking



We can also identify a positive linear correlation between ranking position and loudness. More specifically, the louder the track, the better the chances of gaining a higher weighted ranking.

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

Linear Regression: Energy & Ranking



Looking at the graph above, we can see that there is no relationship between a song's energy and its ranking position. The almost horizontal line connecting the observations indicates that there is no influence on the dependent variable.

Following the demonstration of the existence or absence of linear correlations between these three variables and ranking position, we used a t-test to determine whether these variables, taken individually, influence the top ten song selection in Italy.

```
##
## Welch Two Sample t-test
##
## data: log(danceability * 100 + 1) by TopTen
## t = -1.2209, df = 9.0746, p-value = 0.2529
## alternative hypothesis: true difference in means between group 0 and group 1 is not equal to 0
## 95 percent confidence interval:
## -0.22246120 0.06637812
## sample estimates:
## mean in group 0 mean in group 1
## 4.167942 4.245984
```

A first t-test is carried out to examine the effect of danceability on the population of interest - the top ten songs in Italy. We cannot reject the null hypothesis because the p-value is greater than 0.05. This means that danceability has no effect on the ranking of the first ten songs, which do not have to be danceable in order to be ranked highly in Italy.

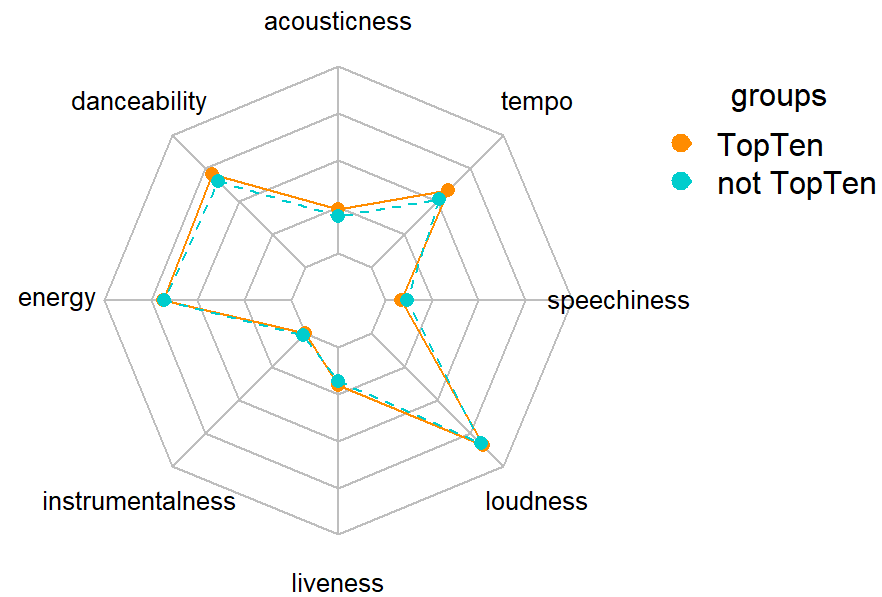
```
##
## Welch Two Sample t-test
##
## data: log(loudness * 100 + 1) by TopTen
## t = -1.6506, df = 9.1678, p-value = 0.1326
## alternative hypothesis: true difference in means between group 0 and group 1 is not equal to 0
## 95 percent confidence interval:
## -0.049179847 0.007620466
## sample estimates:
## mean in group 0 mean in group 1
## 4.42734 4.44812
```

Similarly, due to the high p-value, loudness has no effect on the ranking of the top ten songs in Italy.

```
##  
## Welch Two Sample t-test  
##  
## data: log(energy * 100 + 1) by TopTen  
## t = -0.51931, df = 9.1616, p-value = 0.6159  
## alternative hypothesis: true difference in means between group 0 and group 1 is not equal to 0  
## 95 percent confidence interval:  
## -0.15704695 0.09827626  
## sample estimates:  
## mean in group 0 mean in group 1  
## 4.193915 4.223300
```

We perform the same test, but this time for energy. Because the p-value is greater than 0.05, we accept the null hypothesis and reject the alternative one. The conclusion is that energy has no effect on the ranking of the top ten songs.

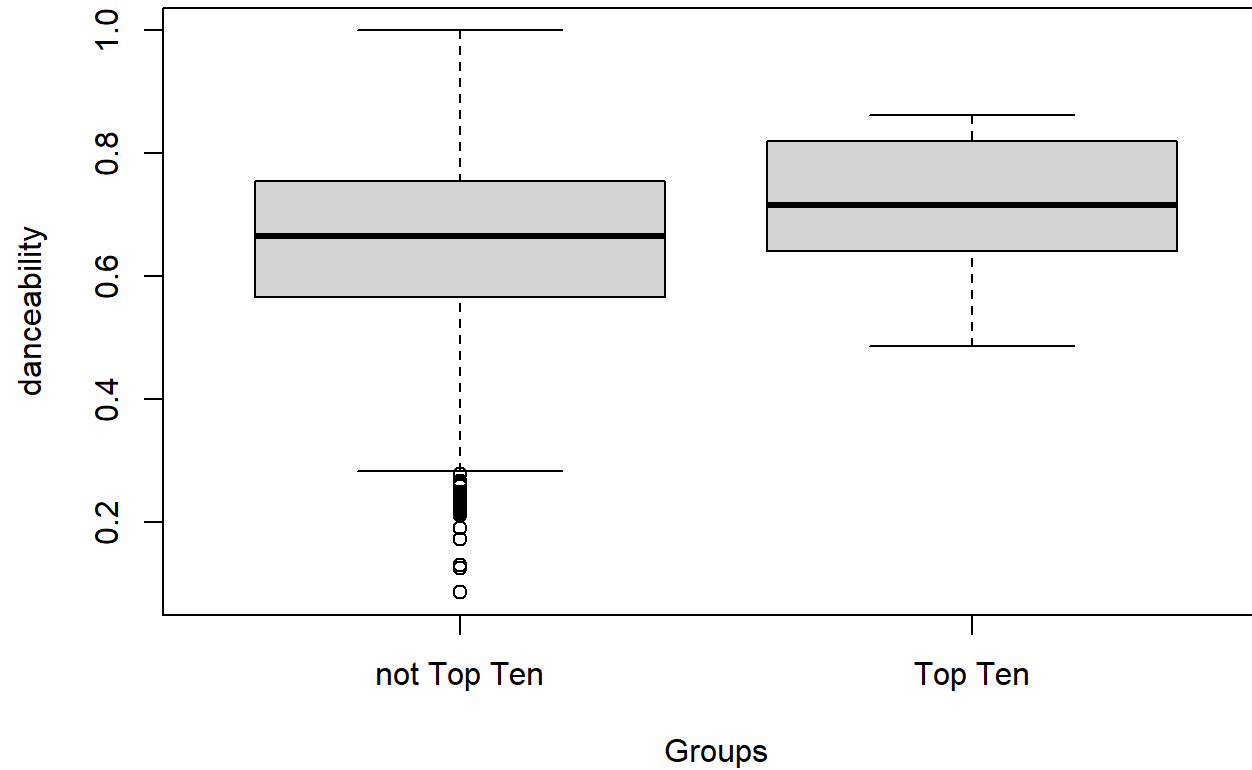
Spiderplot of features in Italy



We ran the radar plot with two groups: Top 10 tracks in Italy and other songs, to get a better visual understanding of the features and their importance to the ranking position of the songs. It is clear that Top 10 songs score higher in terms of danceability, acoustiness, tempo, loudness, and liveness. While scoring lower in all other fields.

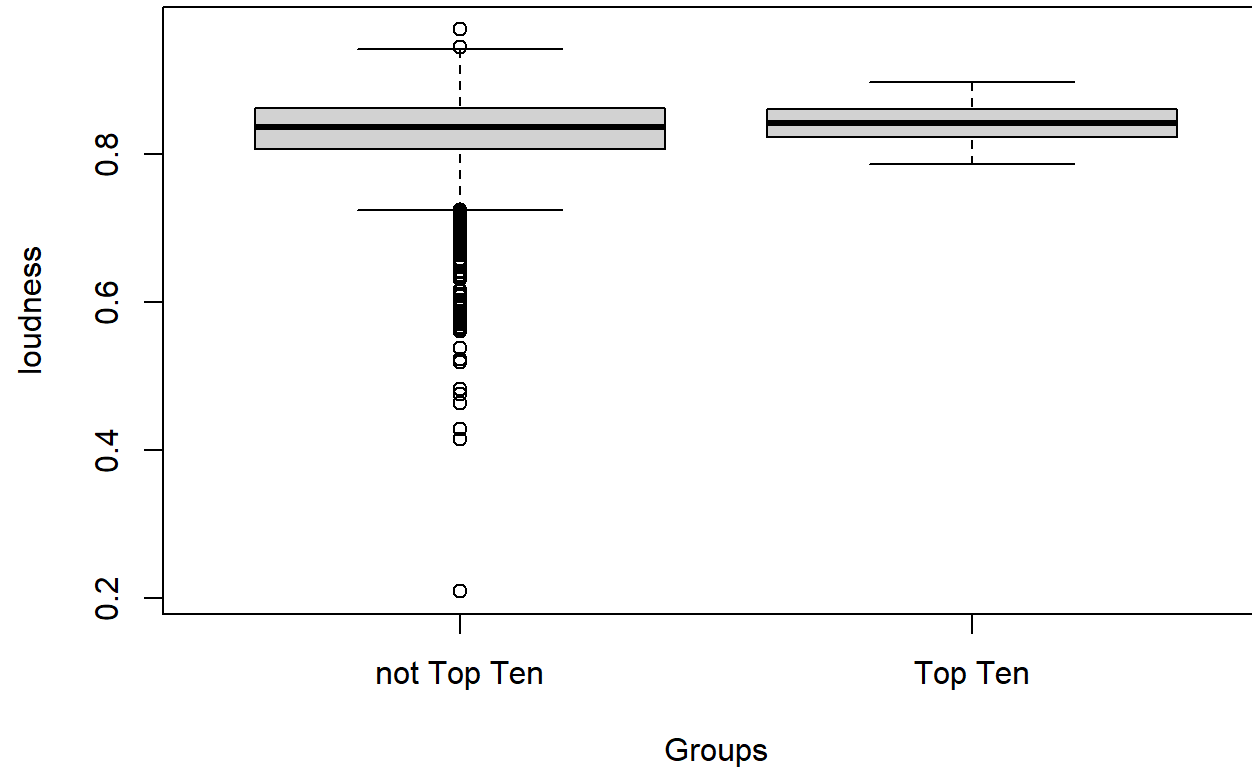
We decided to draw some boxplots comparing top 10 songs vs. non-top 10 songs to further investigate the effect of the three variables.

Boxplot - Danceability



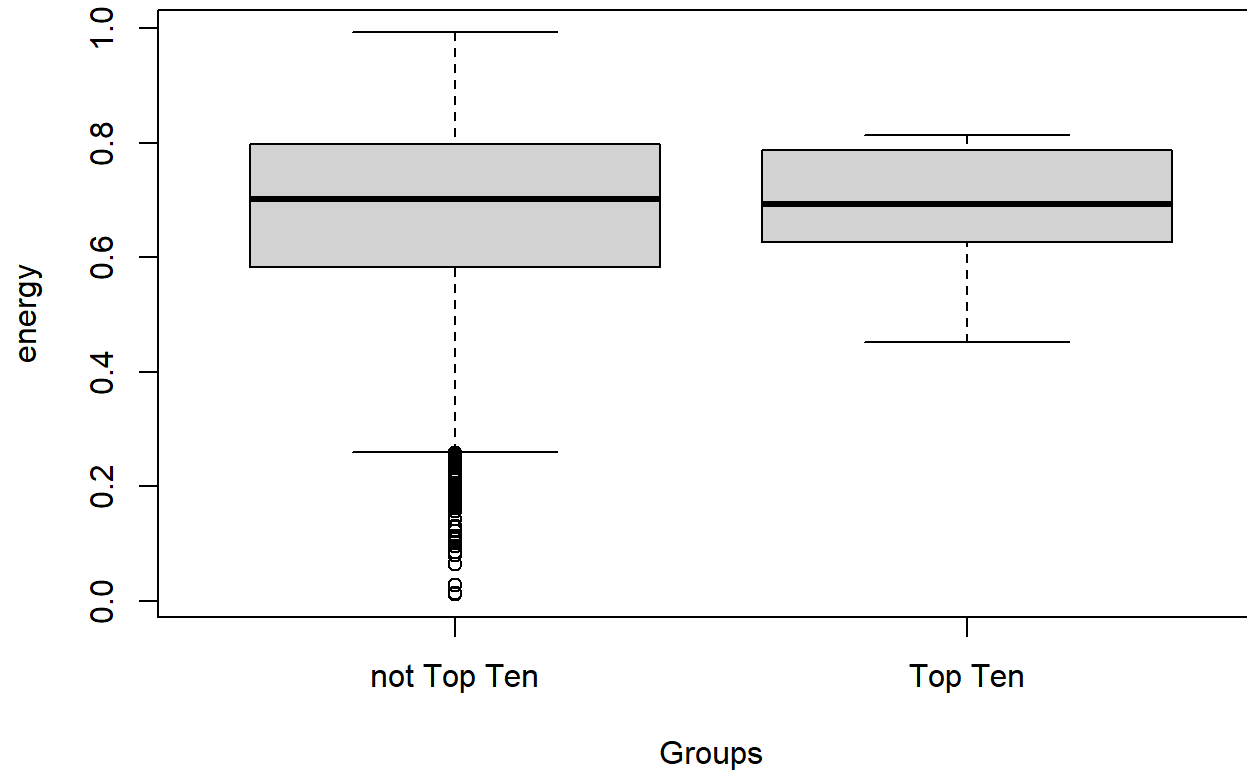
Top 10 songs score higher on average in terms of danceability. Their boxplot shows no outliers and a more compact distribution. The boxplot for non-top 10 songs, on the other hand, shows more dispersion as well as a number of outliers.

Boxplot - Loudness



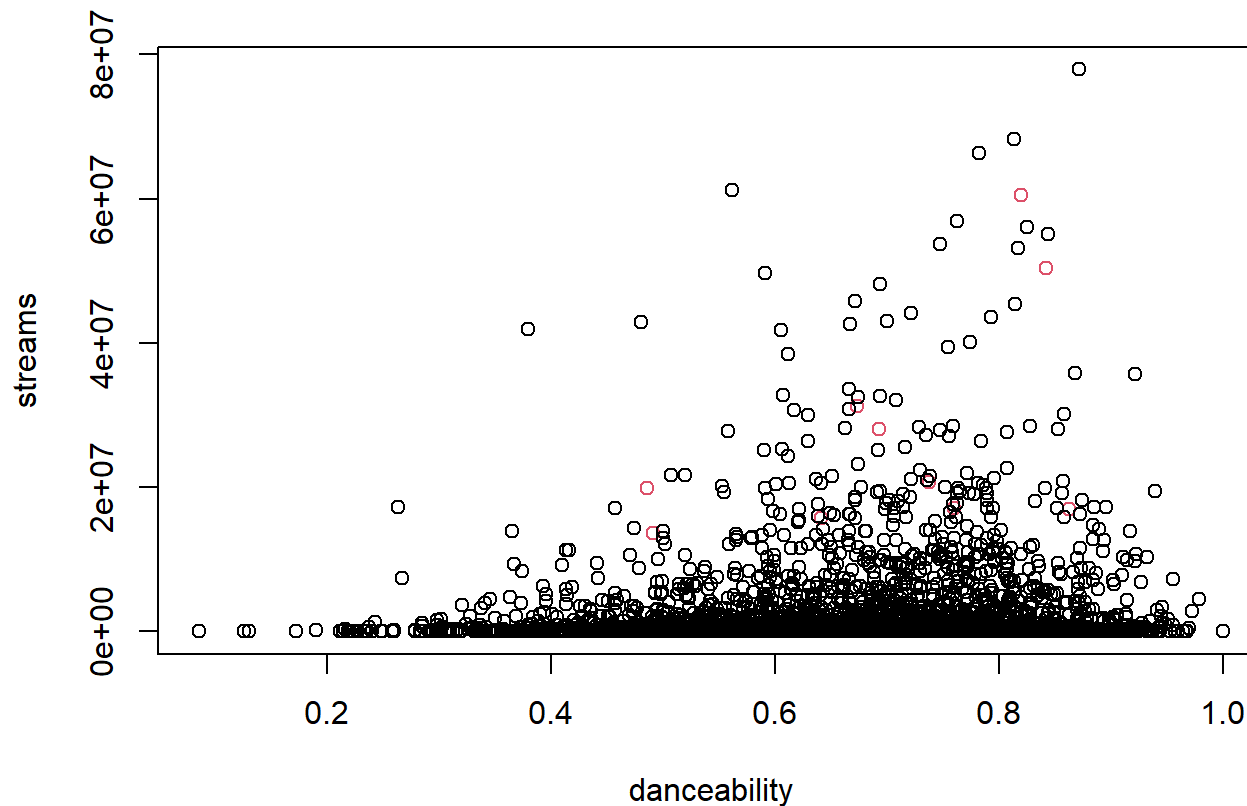
The second boxplot shows that both groups score similarly in terms of loudness, as shown by the medians being very close. Looking at the graph, the non-top 10 songs have a few outliers and a wider dispersion.

Boxplot - Energy



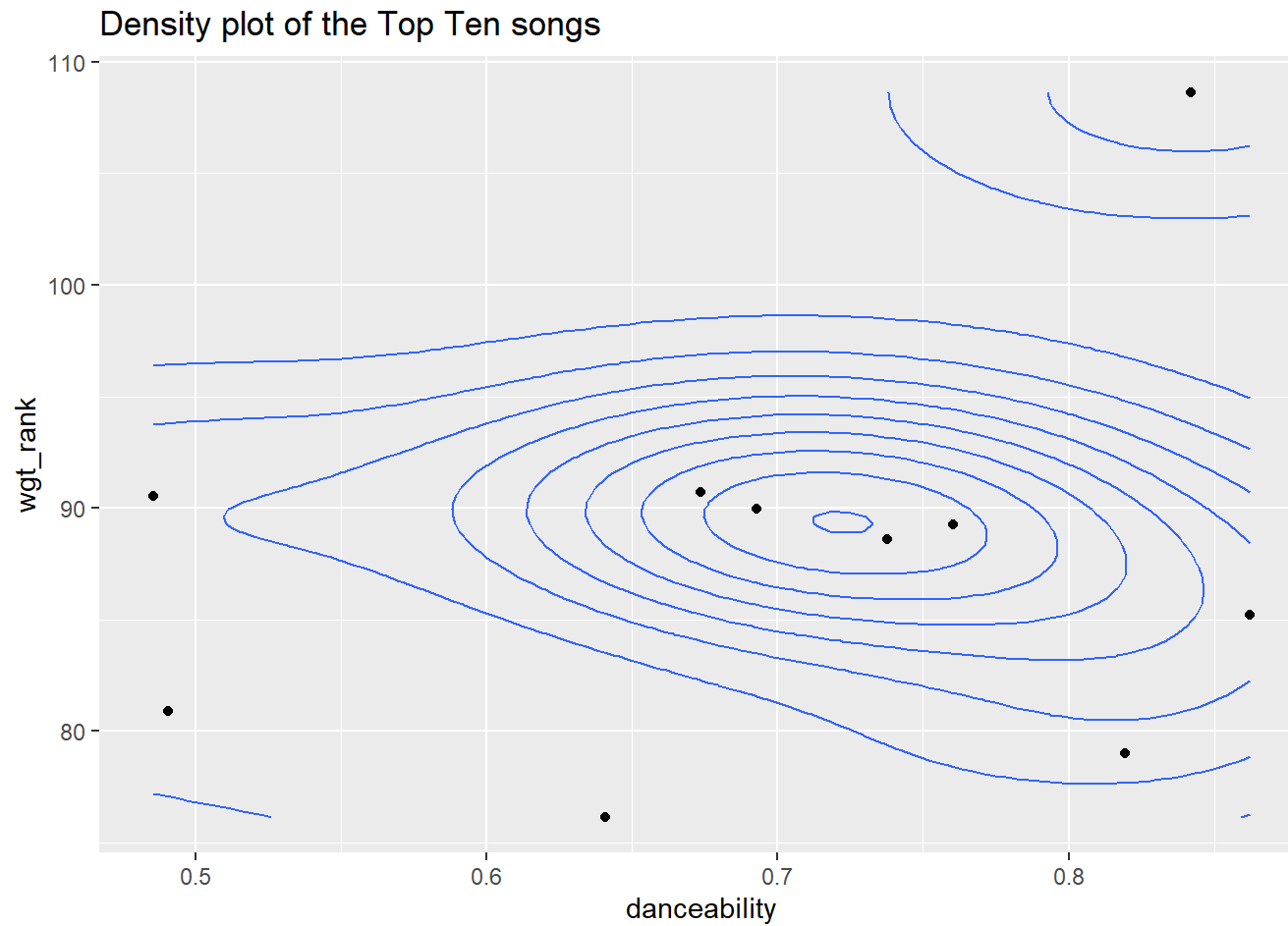
The final boxplot compares the energy levels of top 10 and non-top 10 songs. In terms of level position, the medians of these two groups are similar, yet the mean of the top 10 songs is slightly lower. This indicates that the top 10 songs are generally slightly less energetic, although we know from the analysis before, that this difference is not significant.

Scatterplot: Danceability & Streams



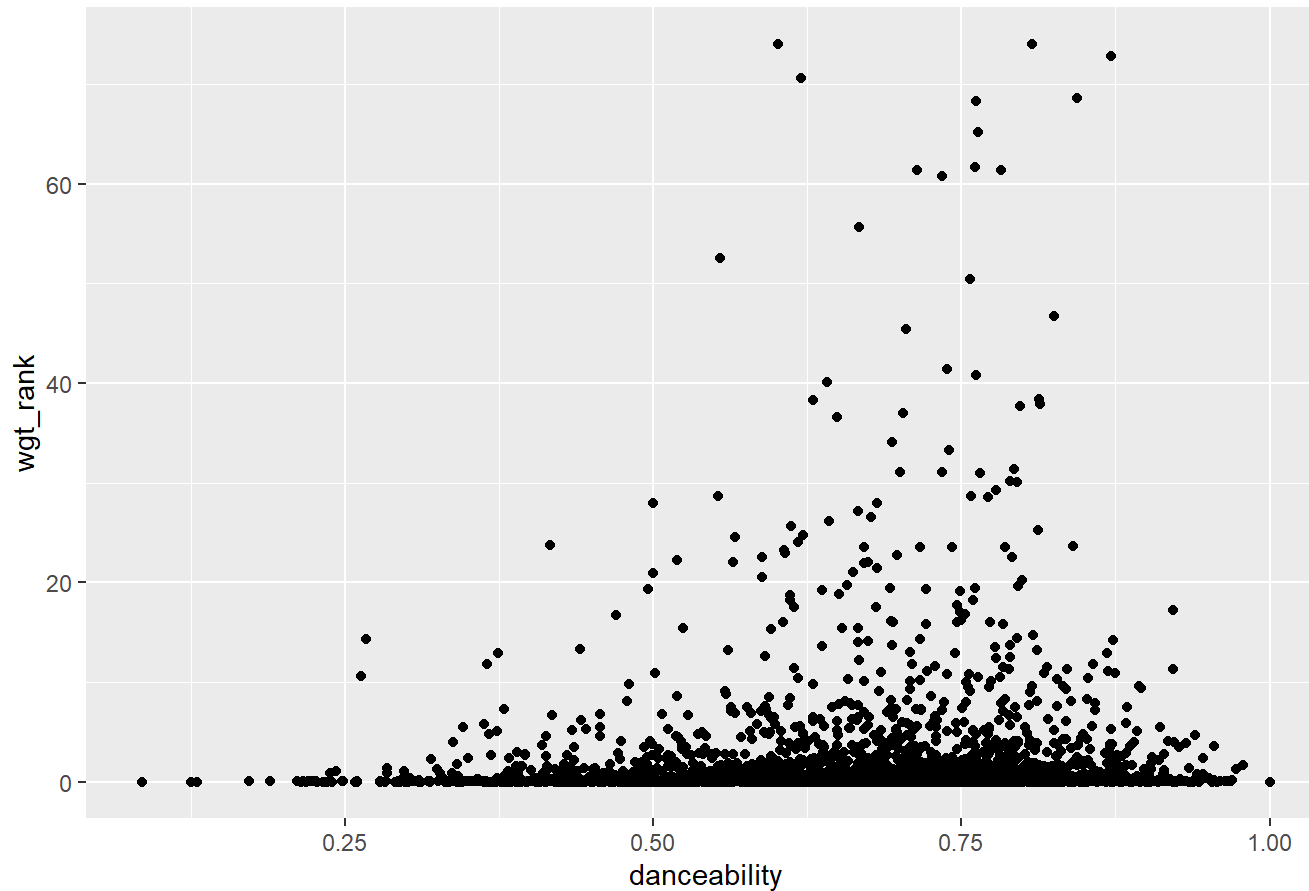
According to the previous analyses, danceability is once again the most important feature in terms of its influence on the weighted ranking. That is why we decided to look into it further also for the case of Italy.

To begin, we created a scatterplot to investigate the relationship between streams and danceability. We can see from this graph that as danceability increases, so does the number of streams. As a result, we can assume that there is a link between danceability and streams.



Second, in order to better understand the relationship between danceability and weighted ranking, we performed a density plot. In the case of Italy, the main density of data is concentrated in the region between 0.65 and 0.85, and there are more contours in this region, indicating that the majority of the data is located there. The data points vary a lot, indicating different danceability scores for the top 10 songs .

Density plot of not Top Ten songs



This final graph examines the relationship between danceability and weighted ranking of the songs which are not top 10 songs. Again the density lines of this plot are not visible, because the data is very dense close to zero. We can see that danceability values surrounding 0.75 are high in the ranking position, indicating that danceability is a highly valued feature in higher ranked songs.

USA

Our last country to analyze are the USA. As occurred in the case of previous countries, we started by identifying the most popular tracks between 2015-2021 in the USA. The results are the following:

##	trackName	artistName	streams	wgt_rank
## 1203	One Dance	Drake feat. WizKid feat. Kyla	438173059	125.82939
## 4837	rockstar	Post Malone feat. 21 Savage	207504231	105.05952
## 4477	Closer	The Chainsmokers feat. Halsey	397617778	105.00656
## 4684	HUMBLE.	Kendrick Lamar	549292830	98.08154
## 1825	God's Plan	Drake	321820367	87.56785
## 3811	Sorry	Justin Bieber	221222661	85.47888

The next step is to understand what makes these songs a hit in USA. Once again we decided to run a logarithmic regression with songs that were released only after 2015.

```
##
## Call:
## lm(formula = log(wgt_rank * 100 + 1) ~ log(danceability * 100 +
##      1) + log(speechiness * 100 + 1) + log(acousticness * 100 +
##      1) + log(energy * 100 + 1) + log(instrumentalness * 100 +
##      1) + log(liveness * 100 + 1) + log(loudness * 100 + 1) +
##      log(tempo * 100 + 1) + log(valence * 100 + 1) + time_signature +
##      mode + key + explicit + duration, data = chart_us_year)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.0229 -1.7276 -0.4832  1.4778  6.7157
##
## Coefficients:
##
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      -5.736e+00  2.835e+00  -2.023  0.043118 *
## log(danceability * 100 + 1)  9.247e-01  1.477e-01   6.263  4.18e-10 ***
## log(speechiness * 100 + 1)  -2.047e-02  5.148e-02  -0.398  0.690948
## log(acousticness * 100 + 1)  1.279e-03  2.936e-02   0.044  0.965265
## log(energy * 100 + 1)      -4.160e-01  1.676e-01  -2.482  0.013093 *
## log(instrumentalness * 100 + 1) -1.885e-01  5.372e-02  -3.508  0.000457 ***
## log(liveness * 100 + 1)     -1.389e-01  5.781e-02  -2.403  0.016289 *
## log(loudness * 100 + 1)     1.409e+00  7.068e-01   1.994  0.046272 *
## log(tempo * 100 + 1)       9.100e-03  1.365e-01   0.067  0.946834
## log(valence * 100 + 1)      8.100e-02  6.333e-02   1.279  0.200960
## time_signature      7.513e-02  1.257e-01   0.598  0.550002
## mode                -1.535e-01  6.747e-02  -2.275  0.022978 *
## key                  6.284e-03  9.035e-03   0.696  0.486763
## explicit             8.140e-02  8.041e-02   1.012  0.311429
## duration            -1.239e-07  6.587e-07  -0.188  0.850805
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.033 on 3940 degrees of freedom
## Multiple R-squared:  0.02747,    Adjusted R-squared:  0.02402
## F-statistic: 7.951 on 14 and 3940 DF,  p-value: < 2.2e-16
```

Overall, the results are significant, as the p-value is less than 0.05. However, according to the multiple R-squared, the model explains only around 2,7% of the data variability.

Looking at the coefficients, we can observe the following:

- Danceability, instrumentalness, energy, liveness, loudness and mode are the only significant variables with a p-value lower than 0.05.
- Among the significant variables we can distinguish between the ones that positively affect the ranking position of the song – which are danceability and loudness – and those who negatively affect it – which are energy, instrumentalness, liveness and mode.

Hence, we can infer that Americans prefer songs that are especially danceable and loud (coefficients equal to 0.9247 vs 1.409). Instead, tracks that are particularly energetic (-0.4160), instrumental(-0.1885) and live (-0.1389) are less likely to become an American hit.

```
##                                feols(log(wgt_ran..
## Dependent Var.:              log(wgt_rank*100+1)
##
## log(danceability x 100+1)      1.041* (0.2384)
## log(speechiness x 100+1)       0.0084 (0.0185)
## log(acousticness x 100+1)      0.0149 (0.0283)
## log(energy x 100+1)            -0.4697* (0.0990)
## log(instrumentalness x 100+1)  -0.1949* (0.0600)
## log(liveness x 100+1)          -0.1409* (0.0332)
## log(loudness x 100+1)          1.713* (0.4264)
## log(tempo x 100+1)             0.0232 (0.0954)
## log(valence x 100+1)           0.0591 (0.0538)
## time_signature                 0.0851 (0.1519)
## mode                           -0.1412 (0.0766)
## key                            0.0058 (0.0099)
## explicit                       0.1108 (0.1532)
## duration                       -5.07e-7 (7.25e-7)
## Fixed-Effects:                -----
## year                           Yes
##
## _____
## S.E.: Clustered                by: year
## Observations                   3,955
## R2                             0.03984
## Within R2                      0.03314
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

We then ran the logarithmic regression with a fixed effect (the release year). The output demonstrates the within R-squared increased to 3.3% of the explained data variation. It proves once again that the release year of a track significantly compromises the results. This output shows that mode does not have an effect on the ranking when accounting for the release year.

```

##                                feols(log(wgt_ran..
## Dependent Var.:              log(wgt_rank*100+1)
##
## log(danceability x 100+1)      1.018* (0.2073)
## log(speechiness x 100+1)       0.0881 (0.0732)
## log(acousticness x 100+1)     -0.1502* (0.0295)
## log(energy x 100+1)           -0.3650 (0.2626)
## log(instrumentalness x 100+1) -0.0180 (0.0275)
## log(liveness x 100+1)         -0.0895 (0.0663)
## log(loudness x 100+1)         2.869* (0.7547)
## log(tempo x 100+1)            -0.0440 (0.1007)
## log(valence x 100+1)          0.0508 (0.1642)
## time_signature                 0.0305 (0.3154)
## mode                          -0.0333 (0.1483)
## key                            0.0015 (0.0265)
## explicit                      -0.1727 (0.2964)
## duration                      -3.74e-8 (3.64e-7)
## Fixed-Effects:                -----
## year                          Yes
## artistName                    Yes
## _____
## S.E.: Clustered                by: year
## Observations                  3,955
## R2                            0.62885
## Within R2                     0.03840
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

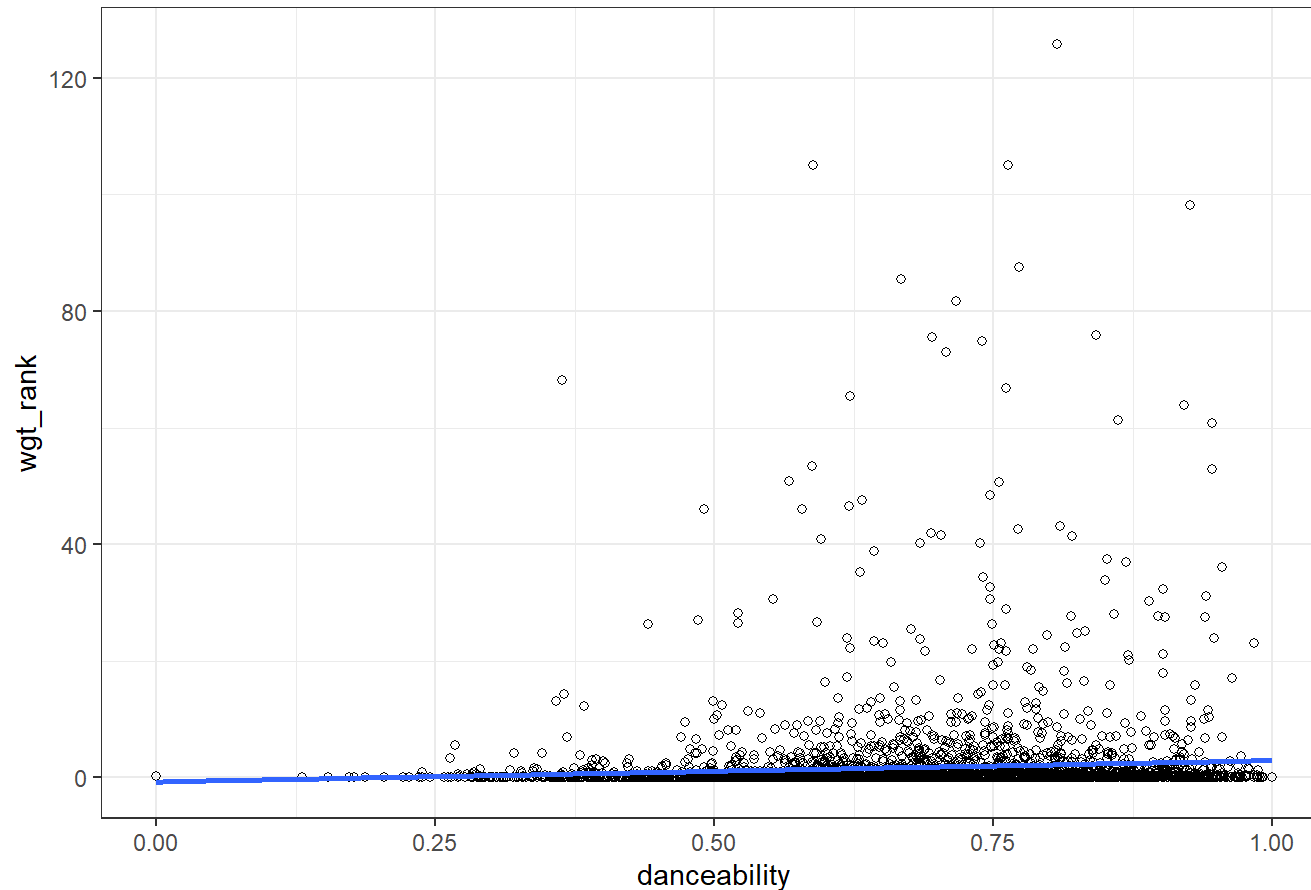
```

The last logarithmic regression that we ran was the logarithmic regression with two fixed effects: the release year and the artist's name. First of all, it is worthy to mention that the multiple R-squared increased again and the model is now able to explain 62.9% of our data. In this model the only significant results are danceability, acousticness and loudness. We can assume that the artist's name affects the relationship between the song's rank and the independent variables. One interesting insight is that in this case the acousticness of a song does indeed influence the weighted ranking in a negative way (-0.1502), while it had a positive effect in the previous regression, and did not have an effect in the other countries. Such a change in sign indicates that there is a different relationship between acousticness and the song's ranking when accounting for the artist.

To sum it up, the results of the American Spotify market are quite similar to the results of the previous countries. As we can see from the regressions, the only difference being acousticness. We still consider the three previously presented features - danceability, loudness and energy - in the upcoming analysis to make it comparable.


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

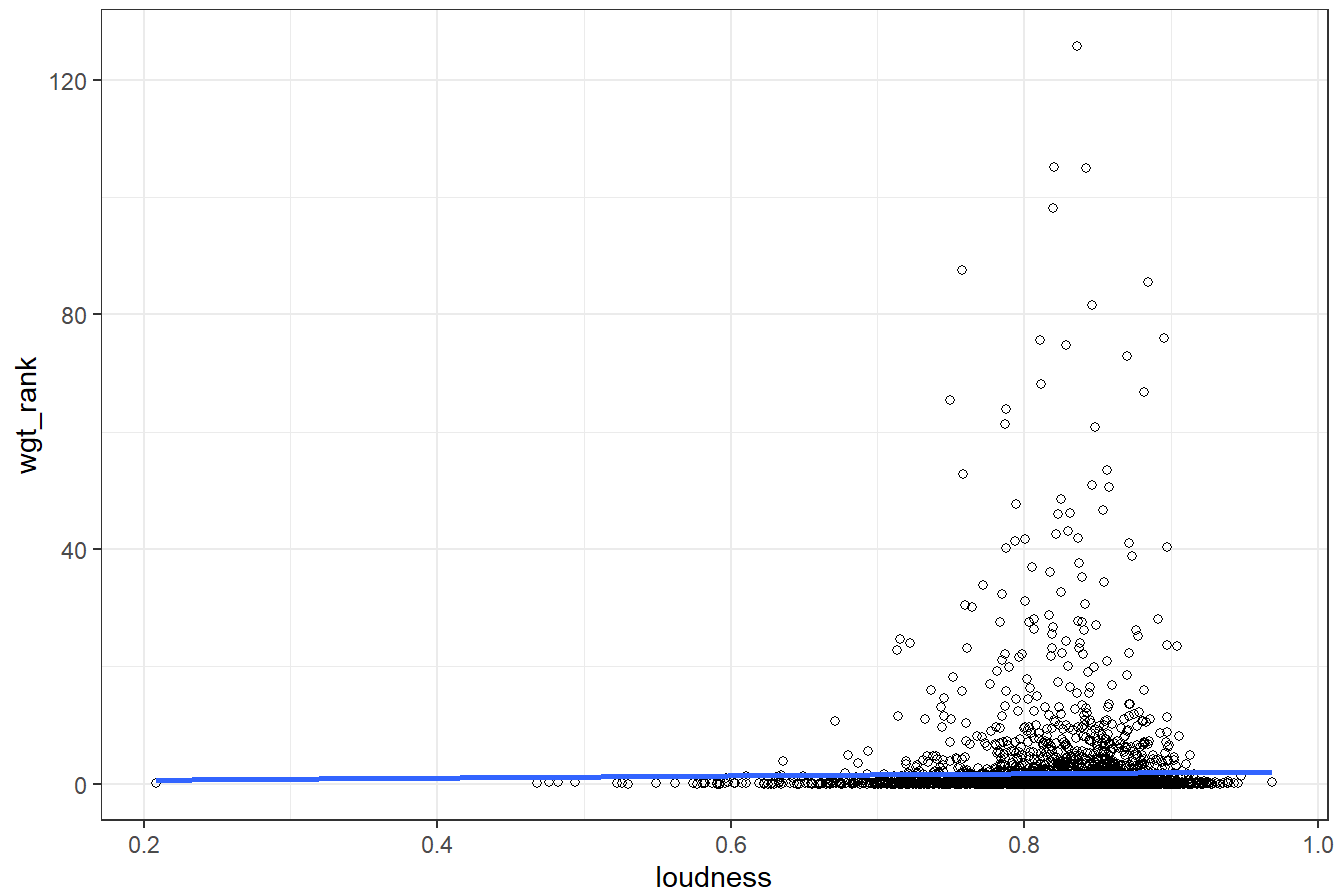
Linear Regression: Danceability & Ranking



Looking at the graphical representation of the relationship between ranking position and danceability, we can identify the positive linear correlation between these two variables. The more a track is suitable for dancing, the stronger is its weighted ranking.

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

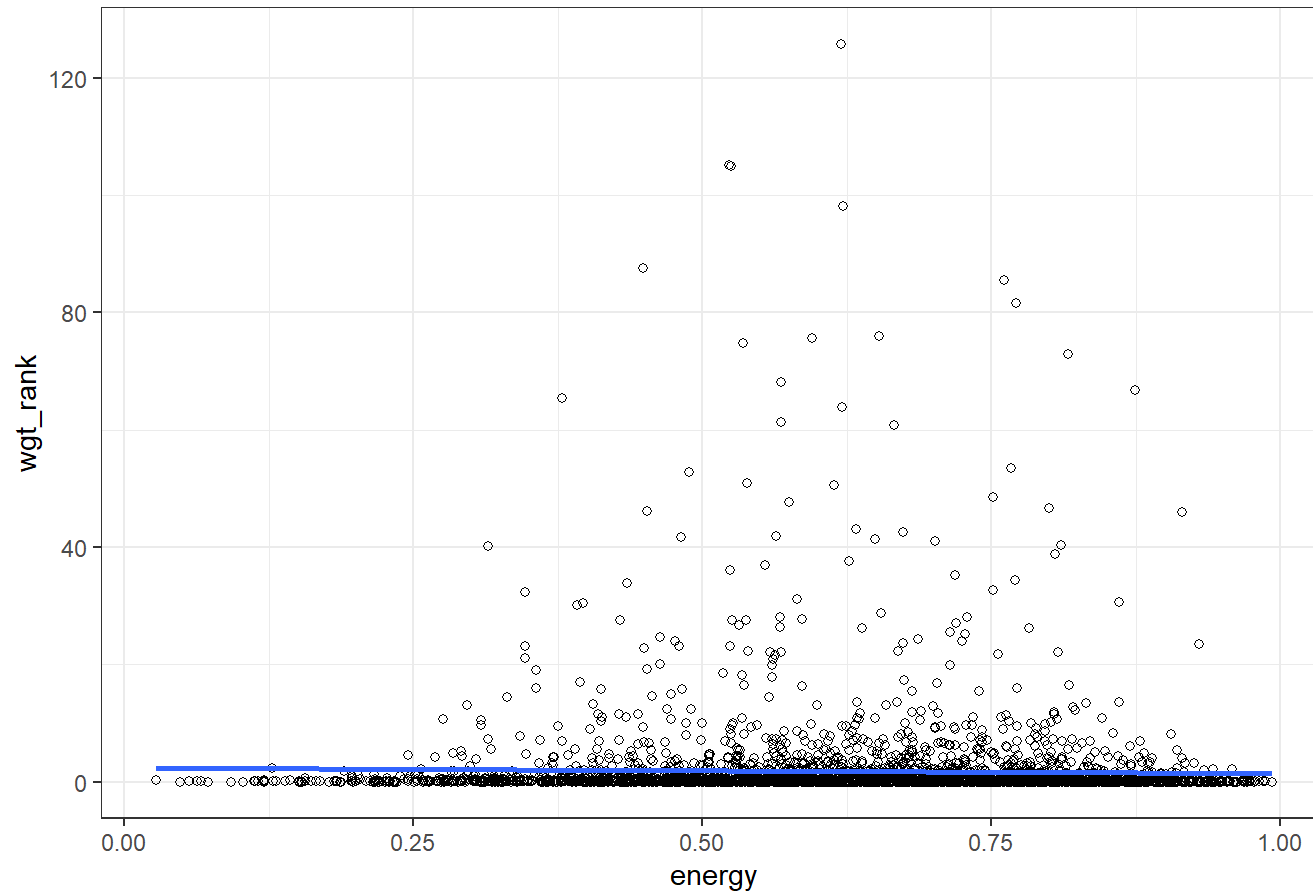
Linear Regression: Loudness & Ranking



Looking at the graph above, we can observe that there is a positive correlation between the loudness of a song and its ranking position.

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

Linear Regression: Energy & Ranking



The graph above shows that there is no linear correlation between the energy of a song and its position in the rank. It shows a slight negative effect, but we know from the analysis above that this is not significant.

After demonstrating the existing or non-existing linear correlations between these three variables – danceability, loudness and energy – and the ranking position, we ran a t-test in order to identify whether these variables taken individually impact the top 10 selection of songs in USA.

```
##
## Welch Two Sample t-test
##
## data: log(danceability * 100 + 1) by TopTen
## t = -3.1053, df = 9.1844, p-value = 0.01231
## alternative hypothesis: true difference in means between group 0 and group 1 is not equal to 0
## 95 percent confidence interval:
## -0.21319832 -0.03380673
## sample estimates:
## mean in group 0 mean in group 1
## 4.202779 4.326282
```

The goal of a first t-test is to analyze the effect of danceability on the population of interest – top 10 songs in the USA. Given that the p-value is lower than 0.05, we can reject the null hypothesis and accept the alternative hypothesis. This means that the difference in means is different from 0, showing that it is very likely that danceability affects the top 10 songs in the US charts. What is more, the sample mean of the top 10 songs is greater than the sample mean of the other songs (4.326 vs 4.202). The conclusion we can get from this is that generally the best 10 songs are more danceable than the other tracks in the chart.

```
##
## Welch Two Sample t-test
##
## data: log(loudness * 100 + 1) by TopTen
## t = -1.0937, df = 9.1129, p-value = 0.3021
## alternative hypothesis: true difference in means between group 0 and group 1 is not equal to 0
## 95 percent confidence interval:
## -0.04871870 0.01692212
## sample estimates:
## mean in group 0 mean in group 1
## 4.419013 4.434911
```

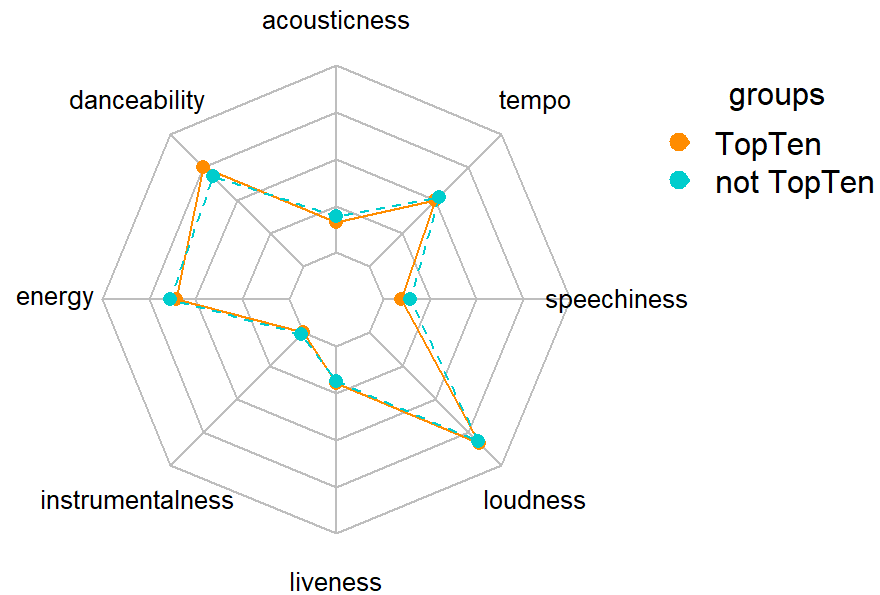
The second t-test analyzes the same for loudness. It shows a p-value of 0.302 and is not significant (>0.05). Therefore loudness has no effect on the top 10 ranking in the US.

```
##  
## Welch Two Sample t-test  
##  
## data: log(energy * 100 + 1) by TopTen  
## t = 0.23686, df = 9.1698, p-value = 0.818  
## alternative hypothesis: true difference in means between group 0 and group 1 is not equal to 0  
## 95 percent confidence interval:  
## -0.1079574 0.1332883  
## sample estimates:  
## mean in group 0 mean in group 1  
## 4.119931 4.107265
```

The third t-test analyzes the effect of energy on the population of interest – top 10 songs in the USA. Given that the p-value is higher than 0.05 (0.818), we accept the null hypothesis and reject the alternative one. The conclusion is that energy does not affect the ranking of the first 10 songs, which therefore do not necessarily need to be energetic in order to be ranked highly.

Next up, we will take a look at the spider graph of the US:

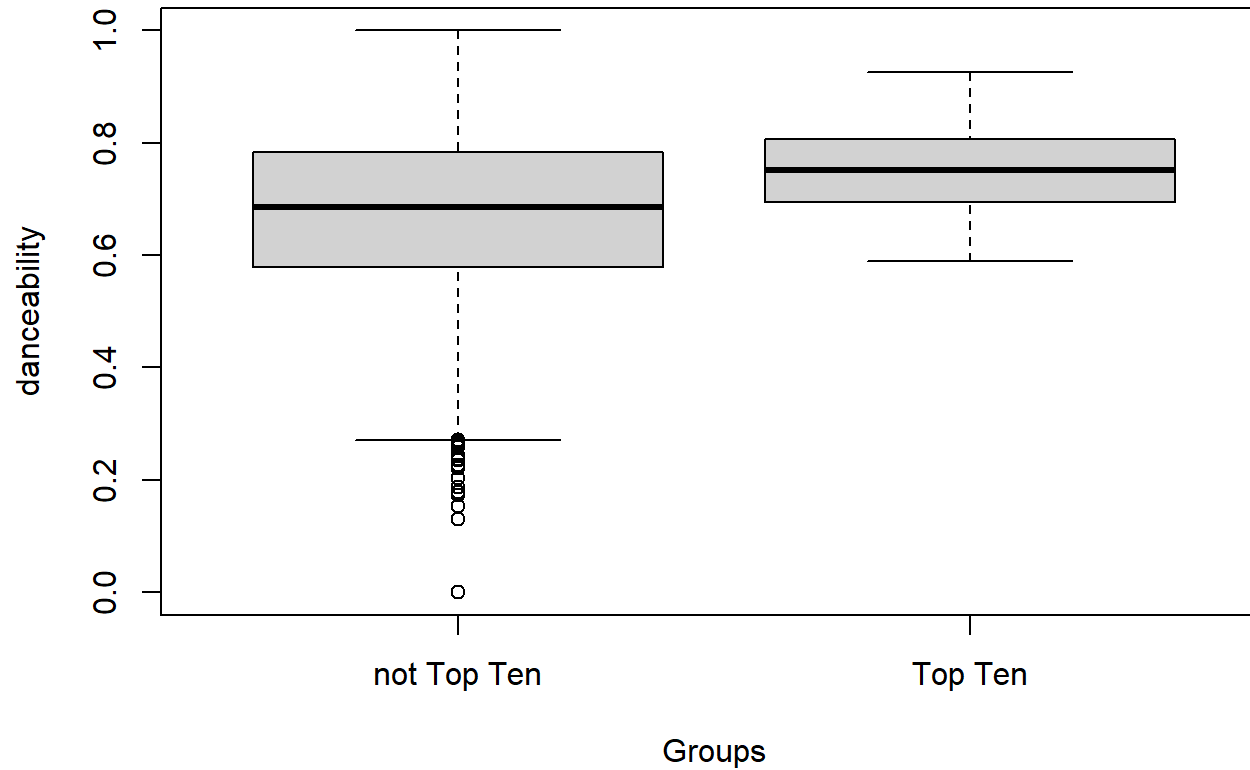
Spiderplot of features in the US



We then constructed the radar plot with our 2 groups: Top 10 tracks in the USA and other songs. It is noticeable that not Top 10 songs seem generally more energetic, faster in terms of pace and score higher in speechiness, while Top 10 songs are more danceable. Moreover, both groups are quite low on instrumentalness, acoustichness, speechiness and liveness, which shows that those features are not characteristic for the US charts.

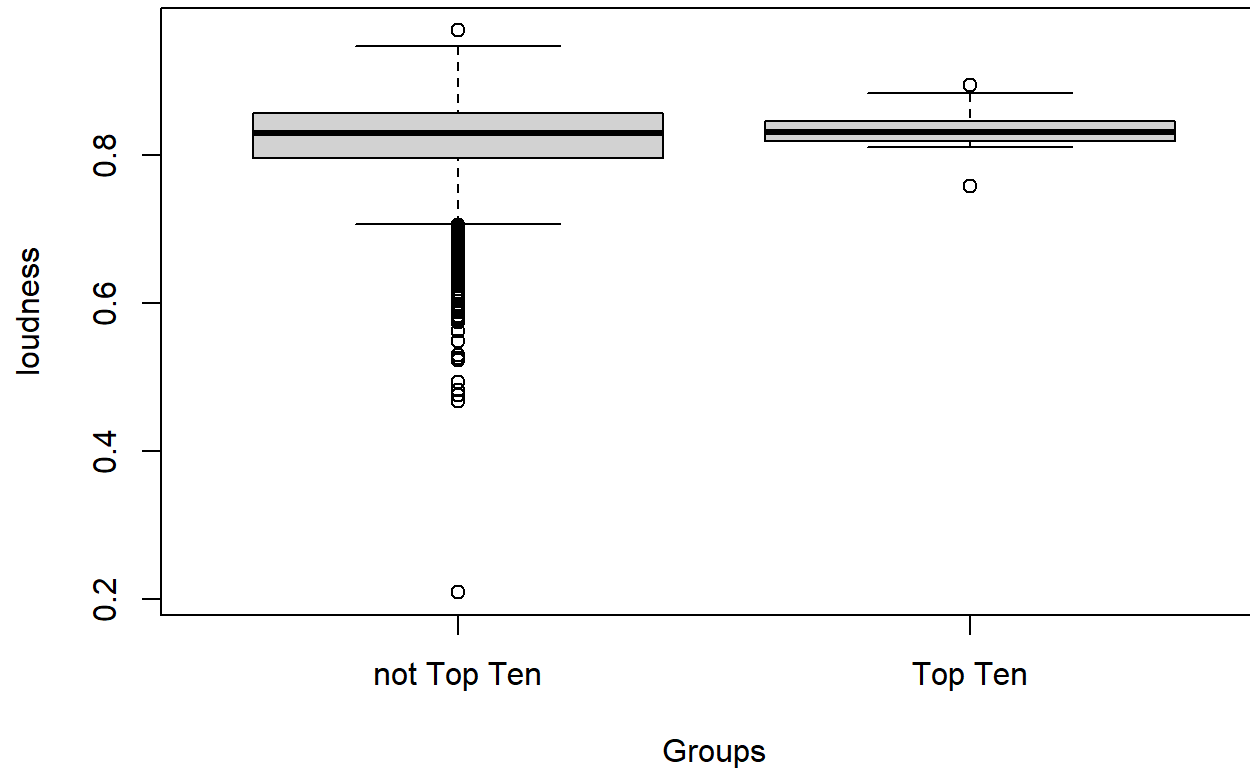
To further investigate the effect of the three variables we focused on – danceability, loudness and energy – we decided to draw some boxplots comparing top 10 songs vs not top 10 songs.

Boxplot - Danceability



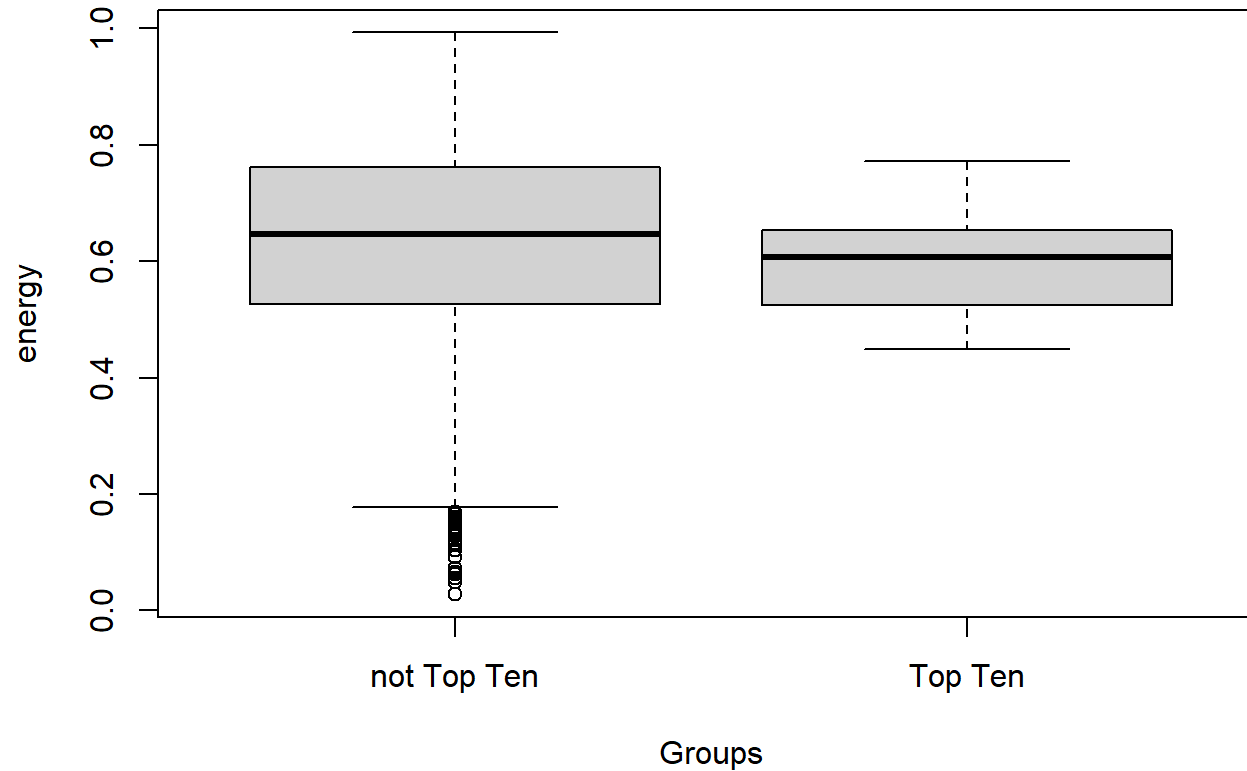
This boxplot shows the difference between top 10 and non top 10 songs in relation to danceability. The medians of the top 10 songs are slightly higher. This tells us that top 10 songs are generally more danceable. The t-test that was conducted before also shows, that this difference is significant.

Boxplot - Loudness



This boxplot shows the difference between top 10 and non top 10 songs in relation to loudness. The medians of these two groups are quite similar. This tells us that top 10 songs are generally not louder than the other ones, but the values of not top 10 songs vary more in loudness, while the top ten songs are more dense in between 0.8 and 0.9.

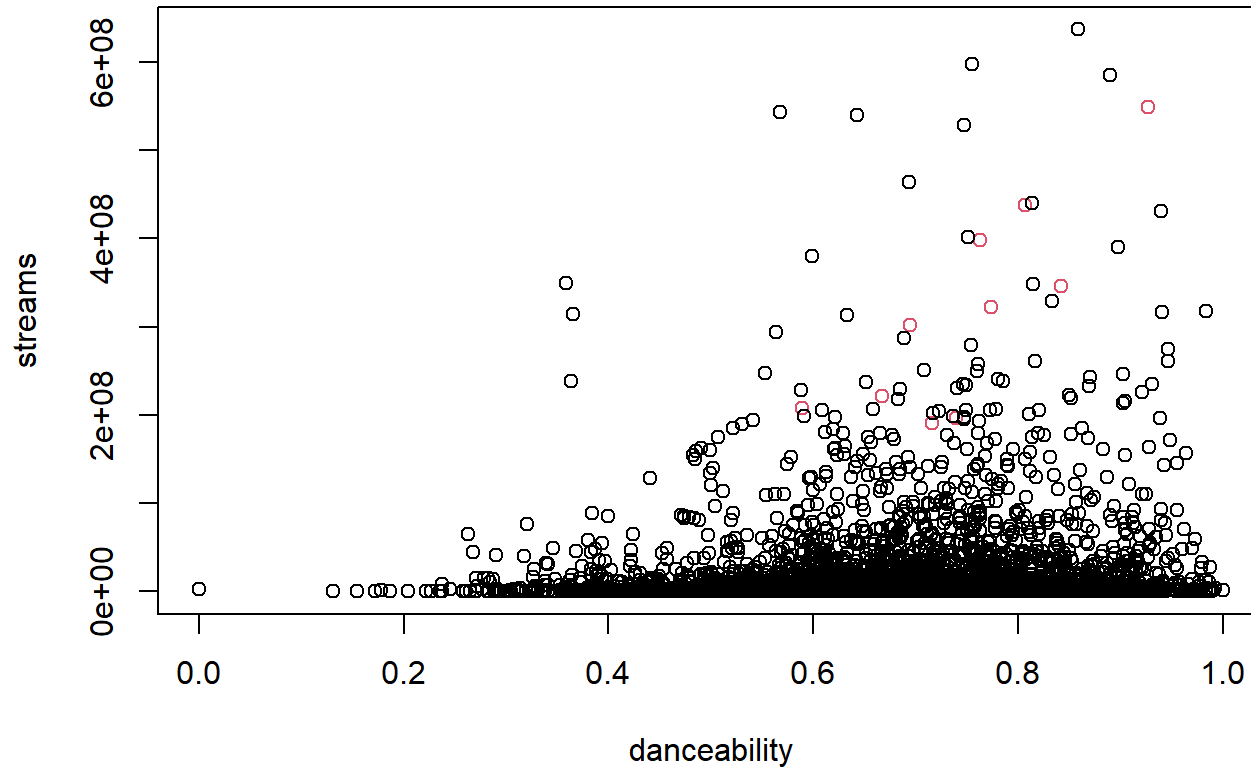
Boxplot - Energy



The last boxplot shows the difference between top 10 and non top 10 songs in relation to energy. The medians of these two groups are somewhat different, the median of the top 10 songs being slightly lower. This tells us that top 10 songs are generally less energetic, although this difference is not significant.

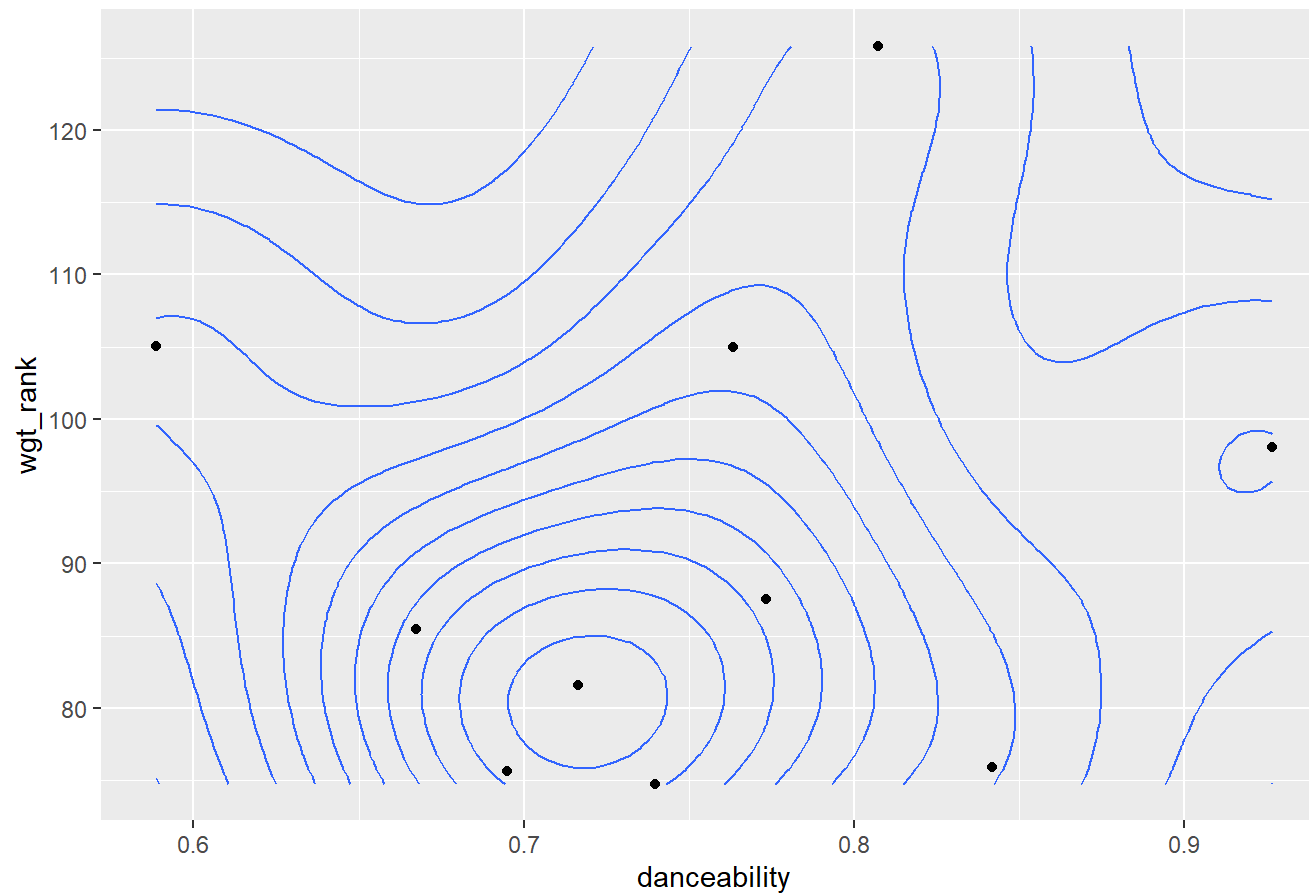
As we can see from the previous analyzes, danceability is once again the strongest feature in terms of its impact on the weighted ranking. That is why we decided to examine it in a more details.

Scatterplot: Danceability & Streams



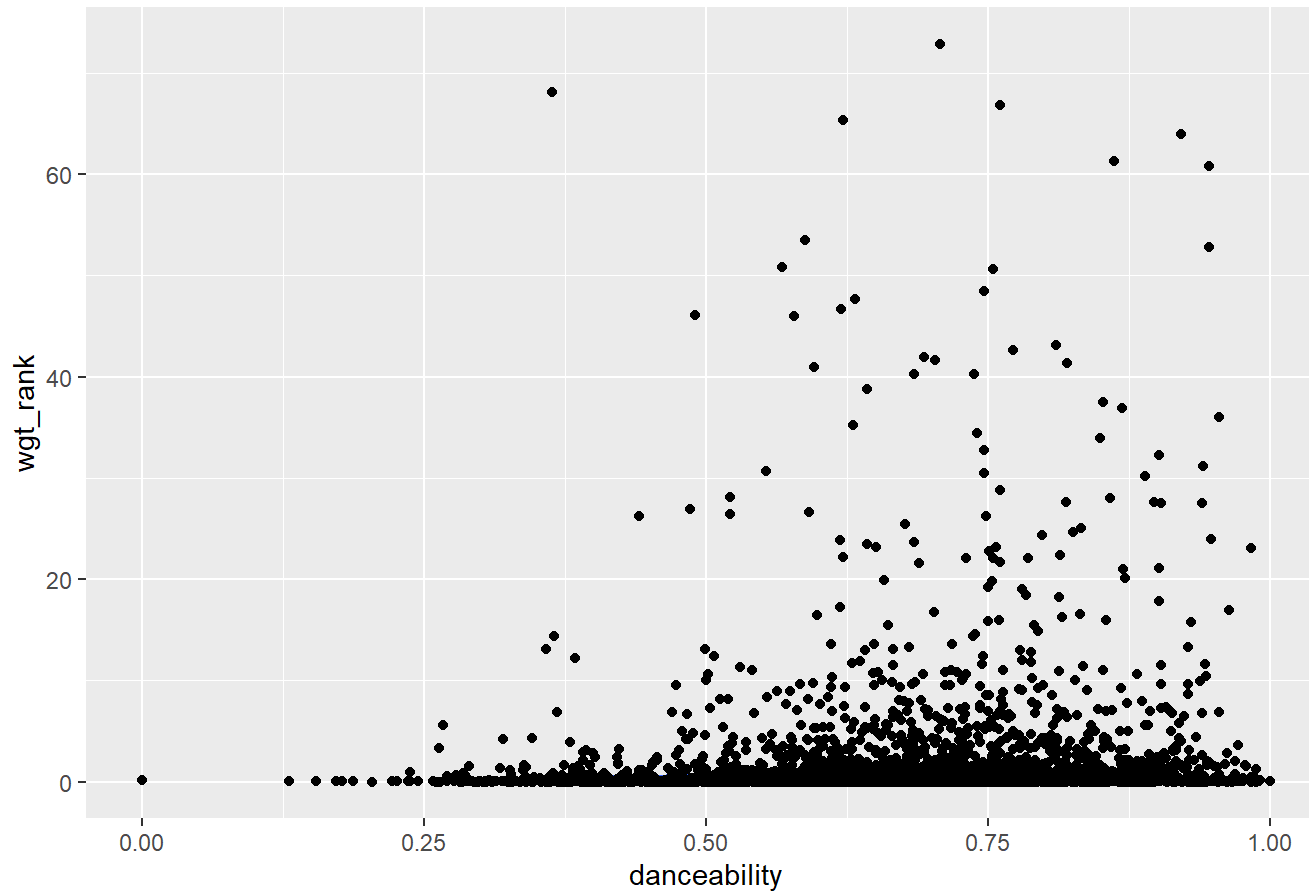
First of all, we computed a scatterplot investigating the relationship between streams and danceability. From this graph we can observe that with the increase in danceability, the number of streams also increases. Therefore, we can assume that there is also an effect between danceability and streams.

Density plot of the Top Ten songs



Secondly, we conducted a density plot in order to better understand the relationship between danceability and the weighted ranking. As we can see, the main density of a data is concentrated in the region between 0.65 and 0.85, but the data varies a lot.

Density plot of not Top Ten songs



The final graph is to look once again at the relationship between danceability and the weighted ranking for the not top 10 songs. We can see that for danceability values surrounding 0.75 the ranking position increases, showing us that danceability is a very valued feature when it comes to higher ranked songs. Again the data is too dense at the bottom to produce a visible density line.

The differences between countries

After presenting the output of our analysis with respect to the three countries individually, we thought it would be interesting to compare them to get a more comprehensive view. We decided to run this comparison analysis on two levels: firstly on the entire dataset, including all songs for the three countries under consideration, and secondly on a reduced dataset, which only includes the top 10 songs.

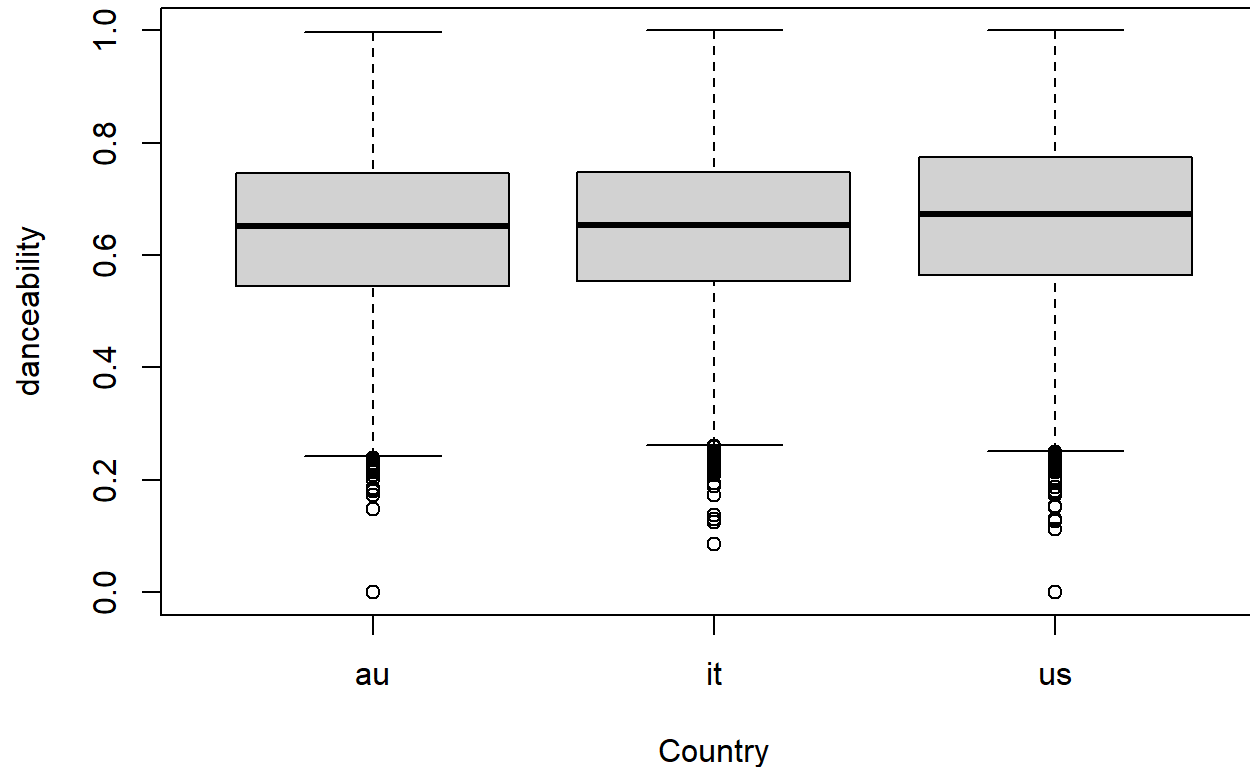
```
##           Df Sum Sq Mean Sq F value Pr(>F)
## region      2   1.74   0.8716   38.54 <2e-16 ***
## Residuals 12930 292.39   0.0226
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
##
## Pairwise comparisons using t tests with pooled SD
##
## data: chart_total$danceability and chart_total$region
##
##      au      it
## it 1      -
## us 1.8e-14 3.8e-12
##
## P value adjustment method: bonferroni
```

In order to perform two statistical tests simultaneously, we decided to execute the Bonferroni t-test, which allowed us to compare different t-test results. We wanted to understand how the three countries behave in relation to the three most interesting variables considered so far – danceability, loudness and energy.

The first test we run was to compare the role of danceability in Australia vs Italy vs US. The test resulted to be significant with p-value lower than 0.05, meaning that we can reject the null hypothesis. Indeed, by looking at the adjusted p-value for the mean difference in danceability scores, we can notice that there are two significant difference: the first one between Australia and US and the second one between Italy and the US.

Boxplot - Danceability by country



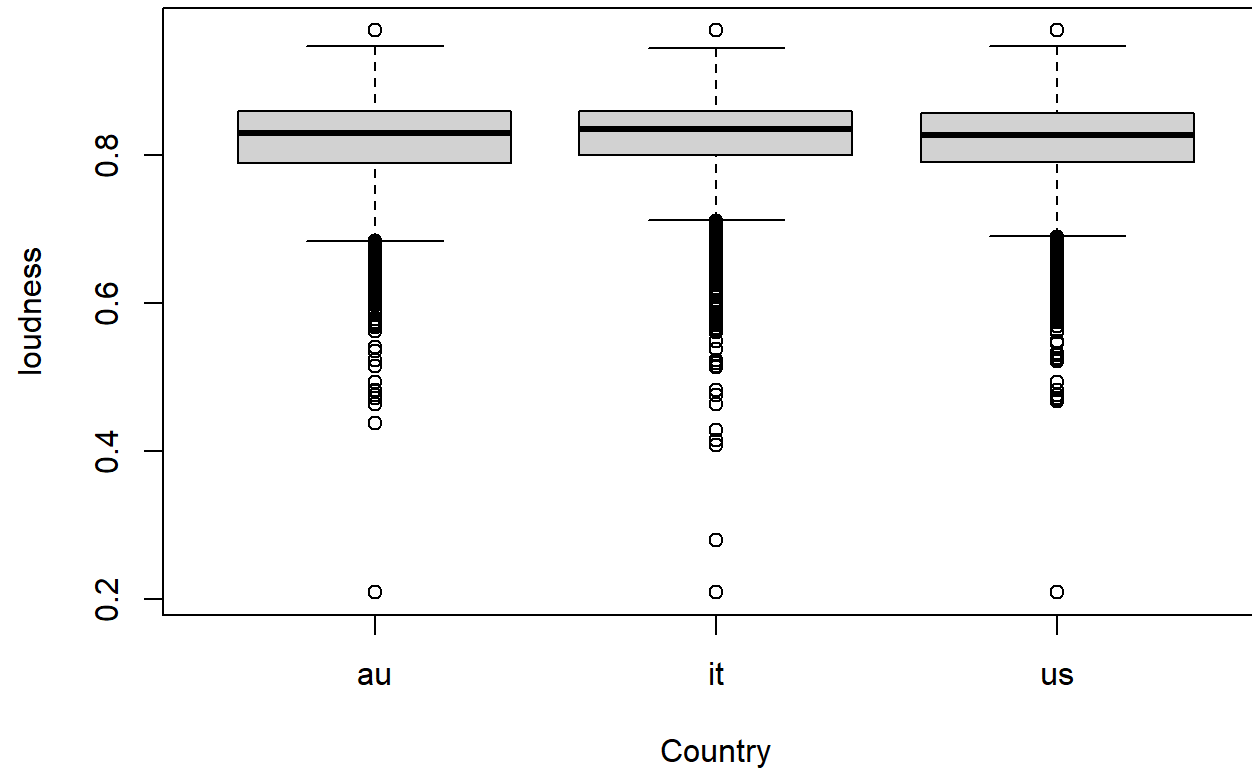
The same results of the Bonferroni t-test can be observed graphically in the following boxplots. Italy and Australia do not present significant differences, as their boxplots are almost identical. On the other hand, comparing US with Australia and Italy it is noticeable that there is a median difference. This result shows us how danceability has a slightly increased importance in the US than in the other two countries.

```
##           Df Sum Sq Mean Sq F value    Pr(>F)
## region      2   0.13  0.06724    18.31 1.15e-08 ***
## Residuals 12930  47.49  0.00367
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
##  
## Pairwise comparisons using t tests with pooled SD  
##  
## data: chart_total$loudness and chart_total$region  
##  
##      au      it  
## it 1.6e-05 -  
## us 0.99    2.5e-08  
##  
## P value adjustment method: bonferroni
```

The following Bonferroni t-test aims at understanding the difference in the role of loudness in the three countries. The p-value is very low, meaning that we can reject the null hypothesis (according to which all means are equal). Looking at the comparisons, we can observe that there is a significant difference between Italy vs Australia and Italy vs the US.

Boxplot - Loudness by country



Following the results of the t-test, we can look at its graphical representation through the use of boxplots. Australia and US look almost identical as their adjusted p-value for the mean difference in loudness scores is equal to 0.99 (almost 1). Instead, Italy shows that loudness is more determinant in making the songs reach the top rankings.

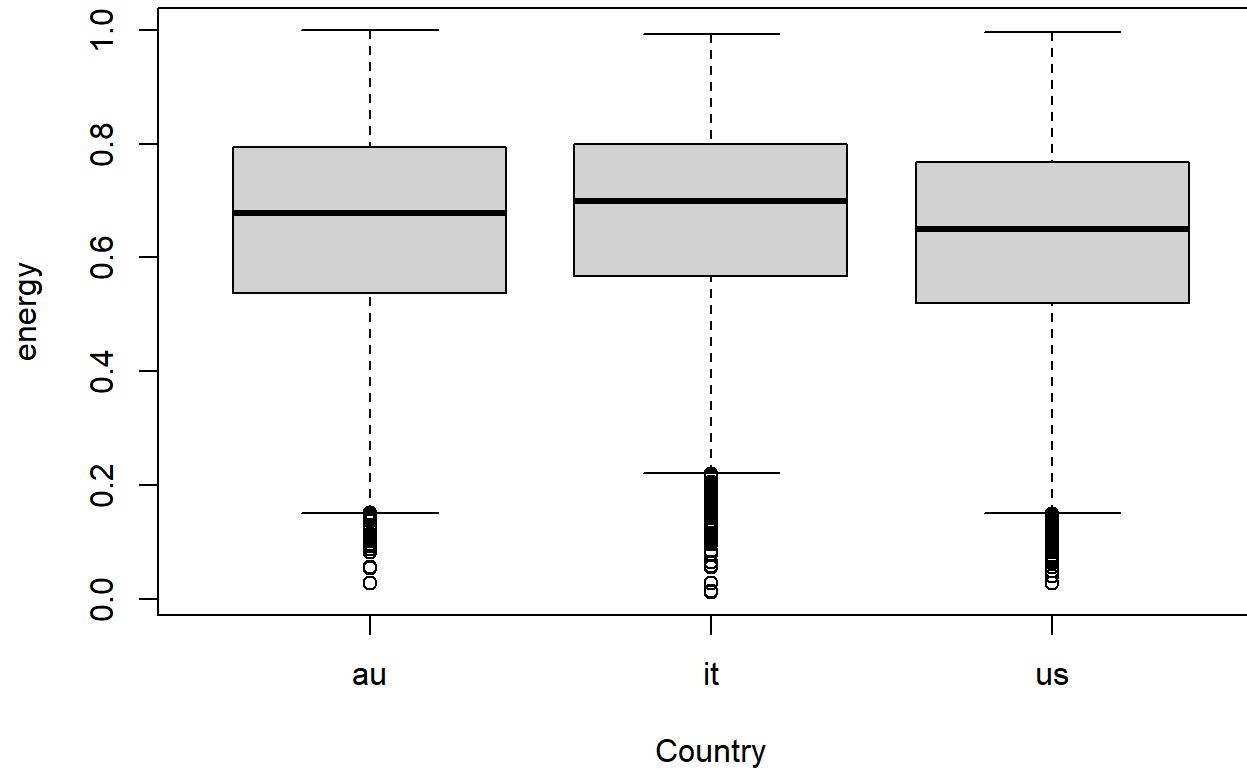
```
##           Df Sum Sq Mean Sq F value Pr(>F)
## region      2    3.3  1.6746   51.58 <2e-16 ***
## Residuals 12930  419.8  0.0325
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```



```
##  
## Pairwise comparisons using t tests with pooled SD  
##  
## data: chart_total$energy and chart_total$region  
##  
##      au      it  
## it 1.0e-05 -  
## us 5.2e-07 < 2e-16  
##  
## P value adjustment method: bonferroni
```

Lastly, we performed the Bonferroni t-test to understand the difference in the role of energy in the three countries. The test is against significant, meaning that we accept the alternative hypothesis assuming that the means are different. By observing the adjusted p-values we can notice significant differences between all countries.

Boxplot - Energy by country

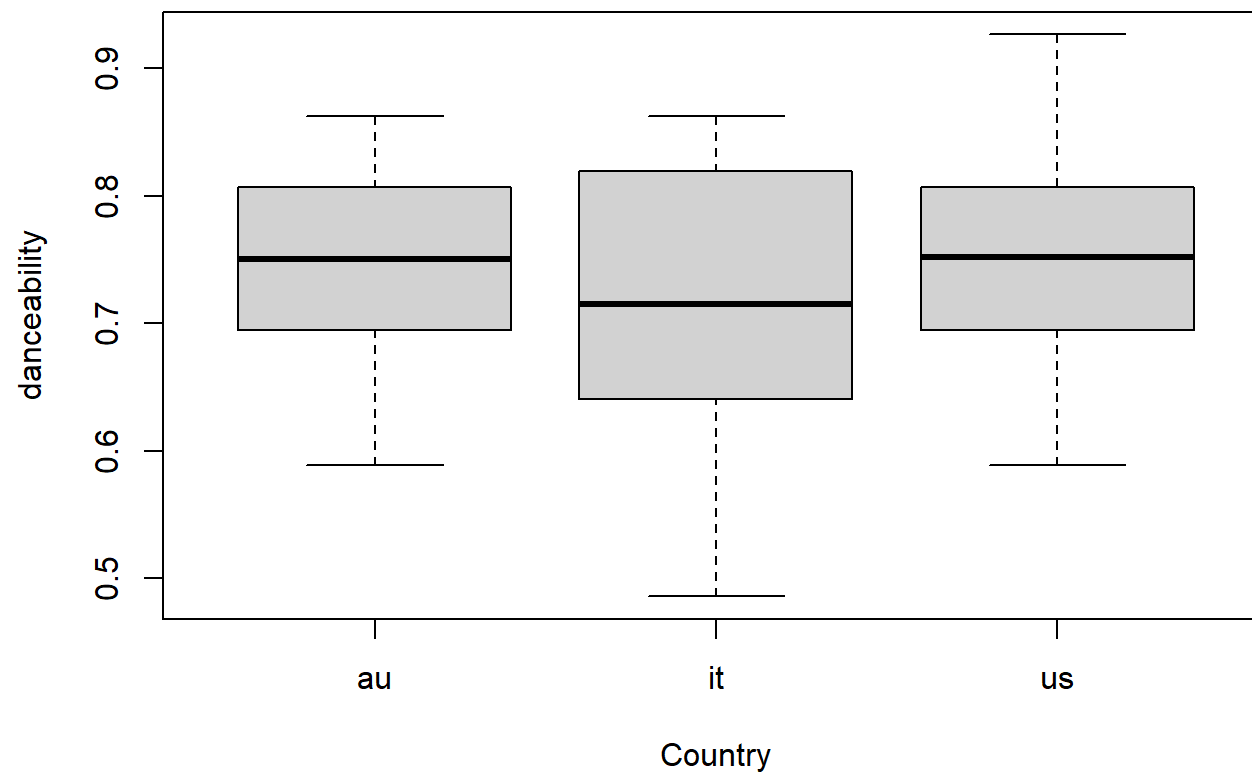


To get a more detailed understanding, we decided to graphically represent the results of the last t-test. We can see that energy is most valued in Italy, while it has the lowest median in US. We can therefore conclude that energy has different values in the three countries, but clearly it is most loved in Italy.

```
##           Df Sum Sq Mean Sq F value Pr(>F)
## region      2  0.01462  0.007312   0.634   0.538
## Residuals  27  0.31133  0.011531
```

```
##  
## Pairwise comparisons using t tests with pooled SD  
##  
## data: chart_short_countries$danceability and chart_short_countries$region  
##  
##    au    it  
## it 1.00 -  
## us 1.00 0.88  
##  
## P value adjustment method: bonferroni
```

Boxplot - Danceability of Top Ten by country

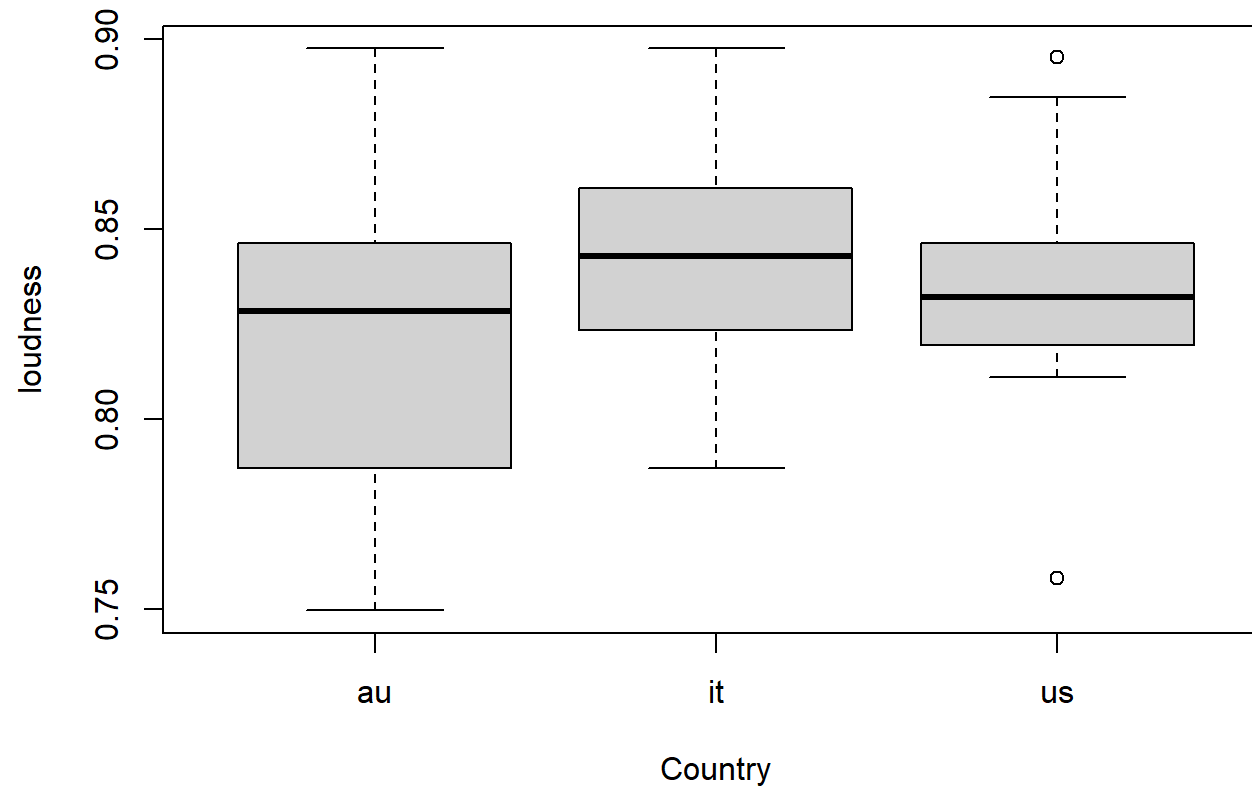


We now conduct the same tests but instead of looking at all the songs we focus on the Top 10. To compare the means and standard deviation of our three matched groups, namely Top 10 songs in Australia, Italy and the US we conduct a pairwise comparisons using t-tests with pooled SD. To evaluate if there is a statistically significant difference between the means of these three independent groups we look at the p-value of the ANOVA. Given the high p-value we accept the null hypothesis that the mean danceability score is the same for each country's Top 10 list of songs.

```
##           Df Sum Sq Mean Sq F value Pr(>F)
## region      2 0.00218 0.001089   0.631   0.54
## Residuals  27 0.04657 0.001725
```

```
##
## Pairwise comparisons using t tests with pooled SD
##
## data: chart_short_countries$loudness and chart_short_countries$region
##
##      au      it
## it 0.81 -
## us 1.00 1.00
##
## P value adjustment method: bonferroni
```

Boxplot - Loudness of Top Ten by country

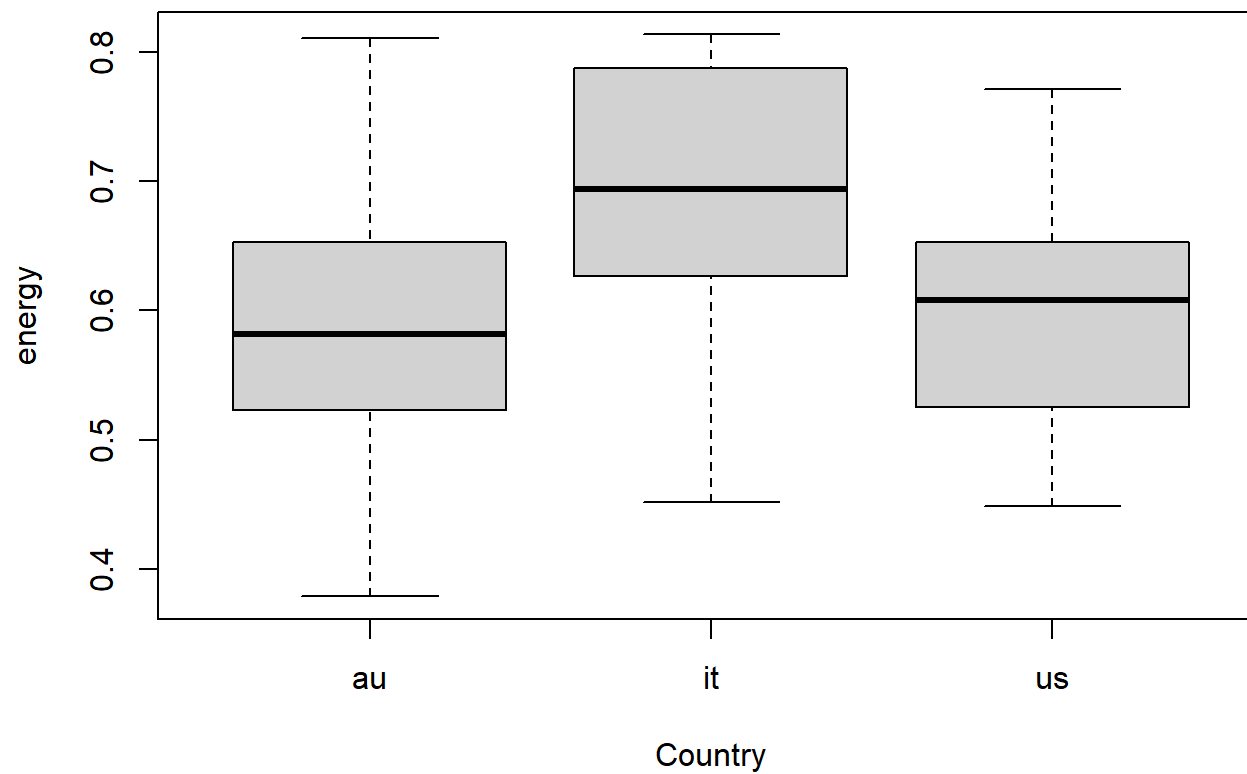


Again, to check whether there is a statistically significant difference between the means of loudness across our three groups, a one-way ANOVA is utilized. Also in this case the high p-value tells us that the mean loudness score is the same for each country.

```
##           Df Sum Sq Mean Sq F value Pr(>F)
## region      2  0.0487  0.02435   1.756   0.192
## Residuals  27  0.3744  0.01387
```

```
##  
## Pairwise comparisons using t tests with pooled SD  
##  
## data: chart_short_countries$energy and chart_short_countries$region  
##  
## au it  
## it 0.27 -  
## us 1.00 0.48  
##  
## P value adjustment method: bonferroni
```

Boxplot - Energy of Top Ten by country

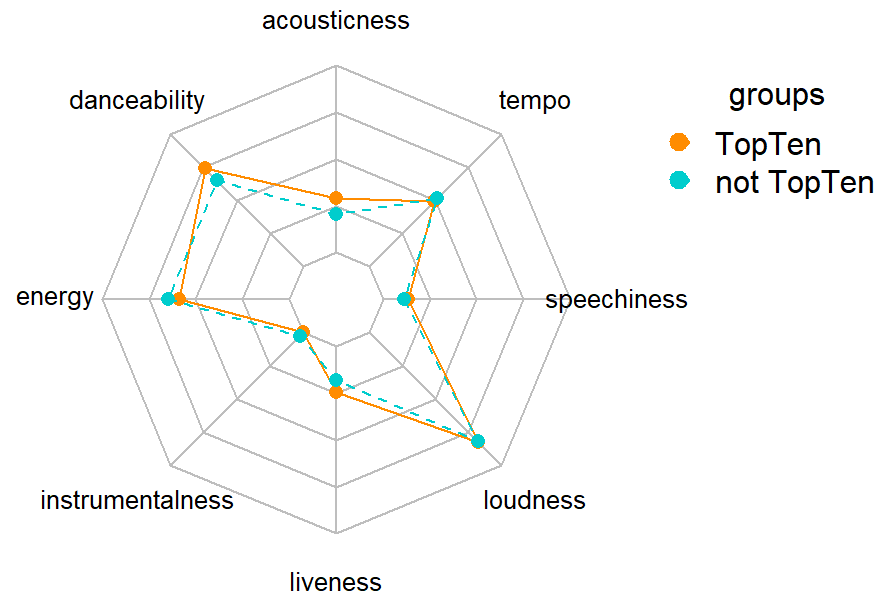


We conduct the same test also to check statistical difference in means of energy. Also in this case, the high p-value tells us that the mean energy score is the same for each country.

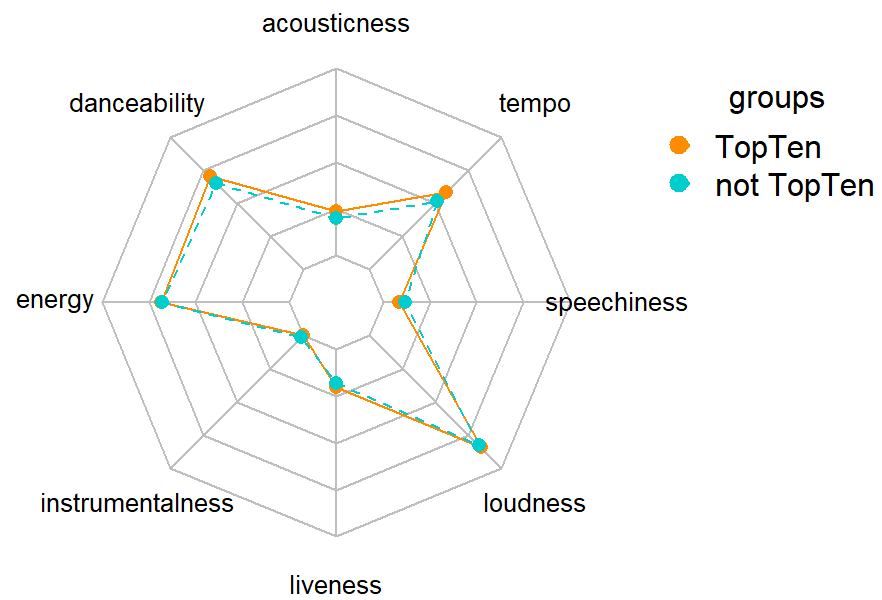
Hence, we can conclude that while these three factors were all significantly different in means when taking into consideration all the songs, they are not while taking into consideration only the Top 10 songs. Hence, all the hit songs are very similar across countries.

The final step in our work is a joint comparison of spider plots for all three countries in order to identify some common and individual patterns.

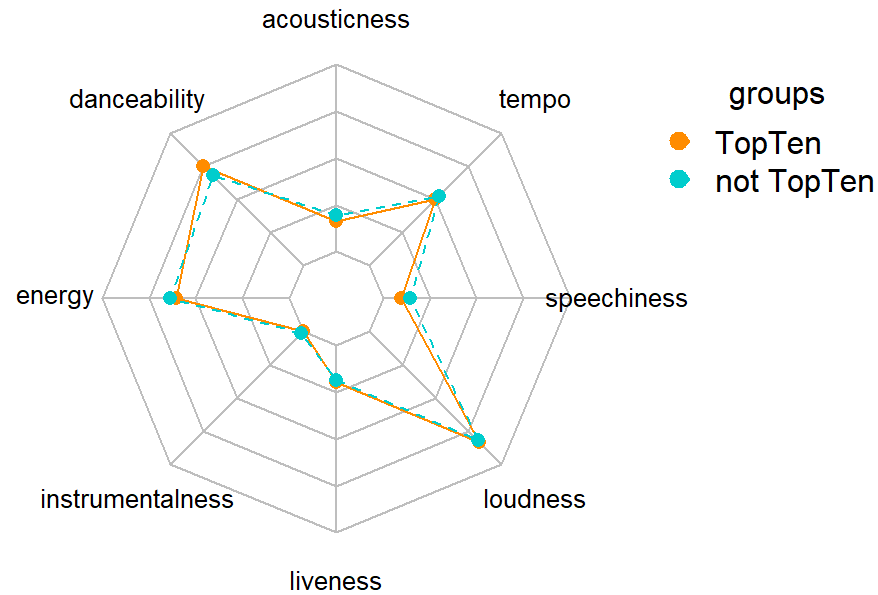
Spiderplot - Australia



Spiderplot - Italy



Spiderplot - USA



Common issues: First of all, we can identify three main, strongest features that are most important and strongly represented in the charts. These are danceability, energy and loudness. What is more, although energy has a strong position, it is more common for a not top 10 tracks. Secondly, there are some features which hold the weakest positions for all three countries and, therefore, are not common or characteristic for songs in the charts of those countries. These are instrumentalness, speechiness, liveness and acoustiness.

Individual patterns: Among the strongest features there are some little individual differences that we would like to draw attention to. The level of danceability and loudness stay the same for all three countries, while the position of energy differs. It is more typical for Italian charts to have more energetic songs (at the same time, both for the top 10 and for non-top 10 tracks). Another thing is the appearance of a tempo feature, that is also more typical for an Italian top 10 chart. Among the weakest features there are also some individual differences. For example, liveness of a song is more inherent for top 10 tracks on the Australian market, while for others (especially for USA) there is not that big difference among top 10 and not top 10. Another interesting pattern occurs with acoustiness. This feature is much more common for top 10 tracks in Australia and Italy, while for the USA this is more inherent for not top 10 songs.

Conclusion

Coming to our conclusion, we should remember our main questions that were stated in this work. We tried to examine what is standing behind the real chart hits and whether there are differences among various countries: Australia, Italy and the USA.

To sum up, our results state that there is a specific type of song that is more likely to achieve the top 10 ranking and that, overall, this type stays the same among different charts throughout the whole world. More danceable and loud songs, but not too energetic, are preferable and will have the greater chance of being liked by the audience. However, every country has its own small individual patterns, that make a track so special for a concrete market. For example, Italian top songs are more energetic in terms of their pace, while for Australian hits danceability is more common than for others.