

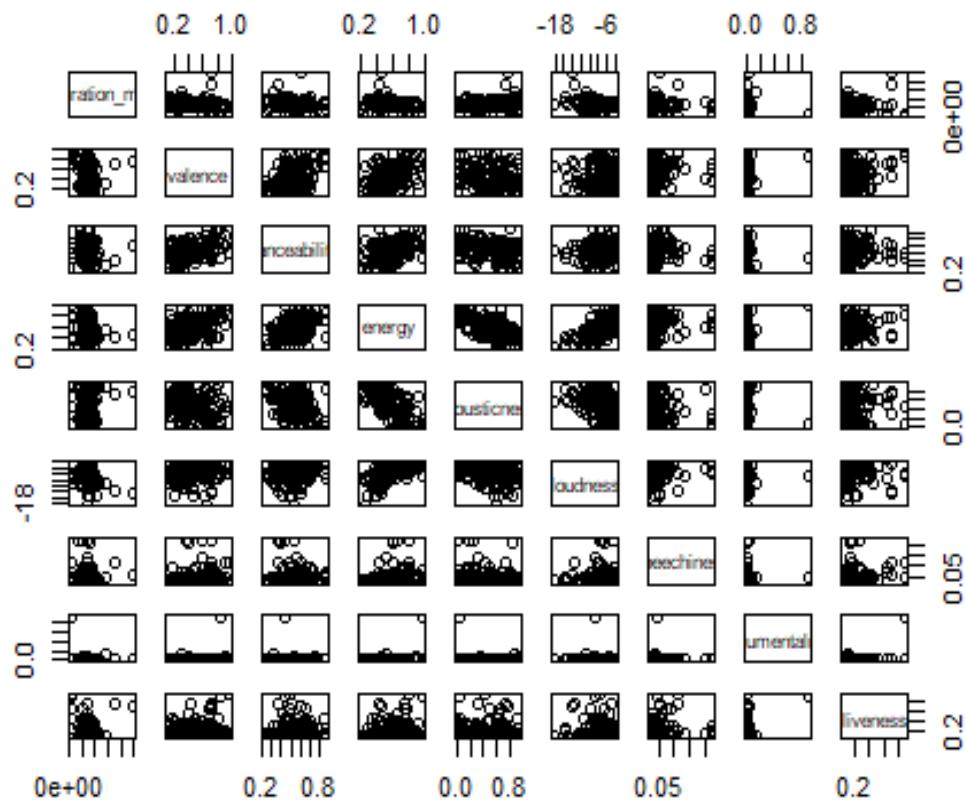
Analysis of Arijit Singh's albums (Spotify Data)

INTRODUCTION:

The artist chosen is a present famous Indian singer, Arijit Singh. He is mostly known for his soothing melodious, soulful and romantic songs. Though he has sung a large variety and genre of songs, his romantic songs are amongst the most famous ones.

EDA:

A pairs plot is plotted amongst the variables listed.

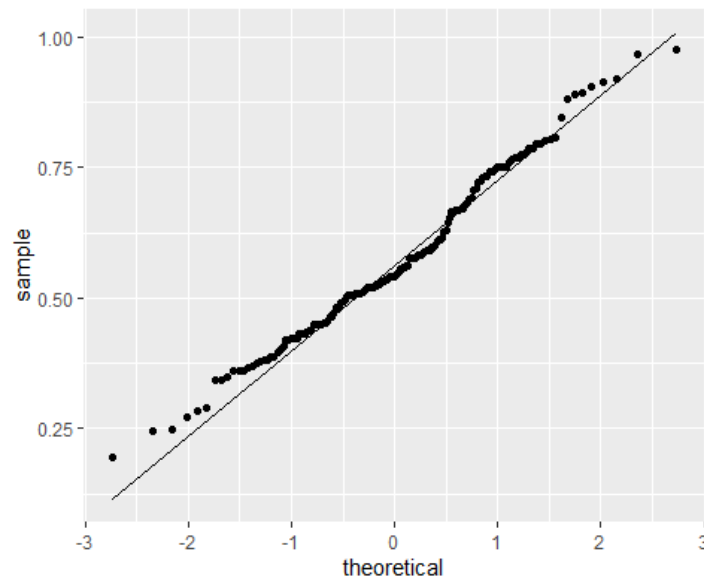


From the above plot, it can be seen that ‘acousticness’ and ‘energy’ are linearly correlated (inversely proportional).

Next, a correlation rank matrix is plotted to confirm the findings from the plot. The pearson's coefficient of -0.65 (between acousticness and energy) is the highest, hence confirming our

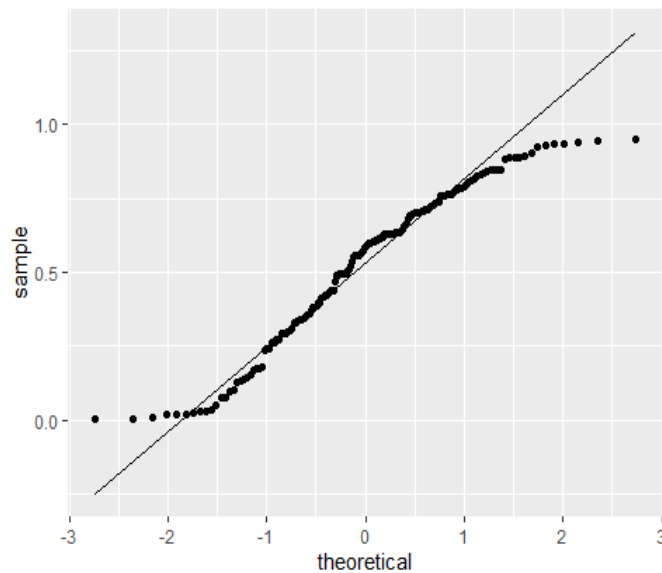
findings. This shows with increase in acousticness there is a decrease in energy of the song. This does make sense as more acoustic the song is its less synthesized by electric means (sounds more like un-plugged version of a track). Hence this may make a more acoustic song as less energetic track.

Normality of energy

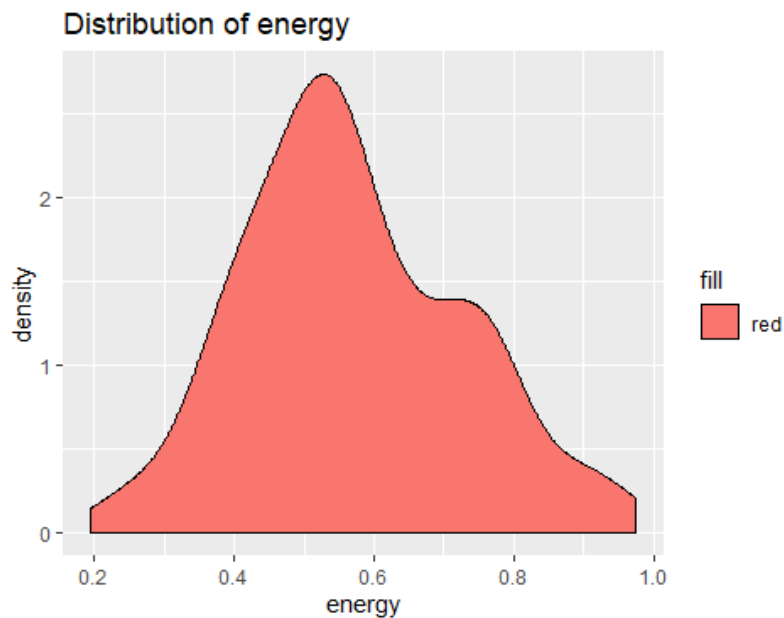


QQ plot of energy is drawn to see that it is quite normal. The distribution deviates from normal at the tails.

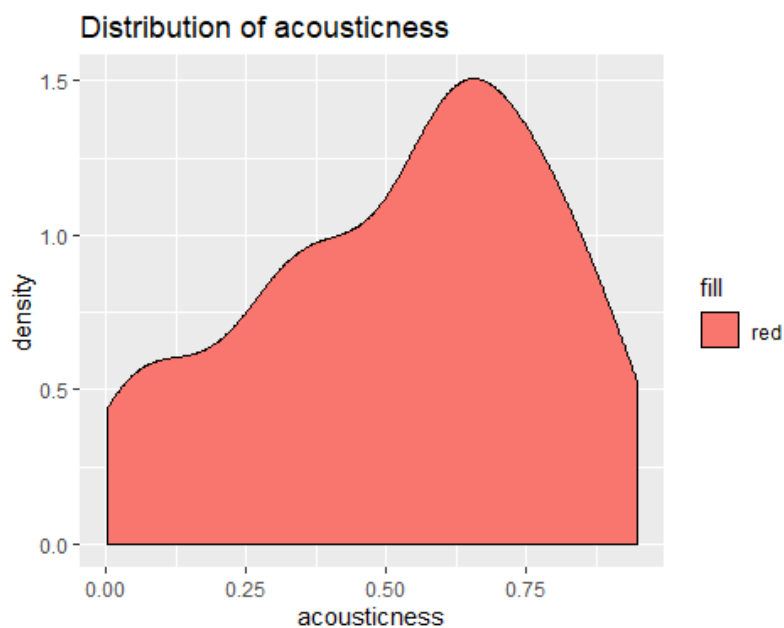
Normality of acousticness



QQ plot of acousticness is drawn to see that it is also quite normal except at the tails it deviates.



The distribution of energy shows that it follows a normal distribution pattern at most parts. This shows that Arijit's most of the songs lie in the mid-energy band.



Acousticness' distribution shows that Arijit's most of the tracks lie on the higher acoustic level (which can be interpreted that most of his songs are not influenced by electric instruments and synthesizers).

Statistical Analysis:

- The 99% confidence interval of **energy** is:

```
## [1] 0.5367533 0.6016467
```

The means of energy level of songs will lie in the above interval with 99% confidence. It further bolsters the statement in previous section that most of Arijit's tracks lie in mid energy level. This goes in accordance with the fact that most of Arijit Singh's songs are melodious than party numbers (assuming party songs will have higher energy levels).

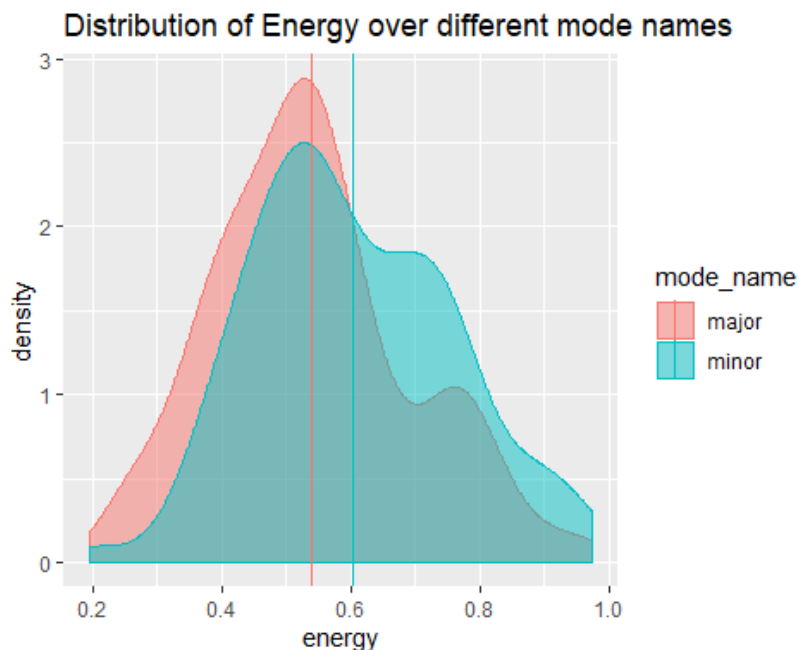
The 99% confidence interval of **acousticness** is:

```
0.4718371 0.5791281
```

The means of acousticness level of songs will lie in the above interval with 99% confidence. This can show that most of his songs are mid-level acoustic.

- 2 sample Hypothesis test on Energy**

Energy variable is chosen to do a pairwise hypothesis test over mode-name of the track. Before performing the test, eyeballing the distribution of **Energy** over the major and minor modes of track.



From the above distribution it can be seen that the distribution of energy level of tracks is not much governed by the mode played (major or minor) on the track, since both the distributions are quite overlapping.

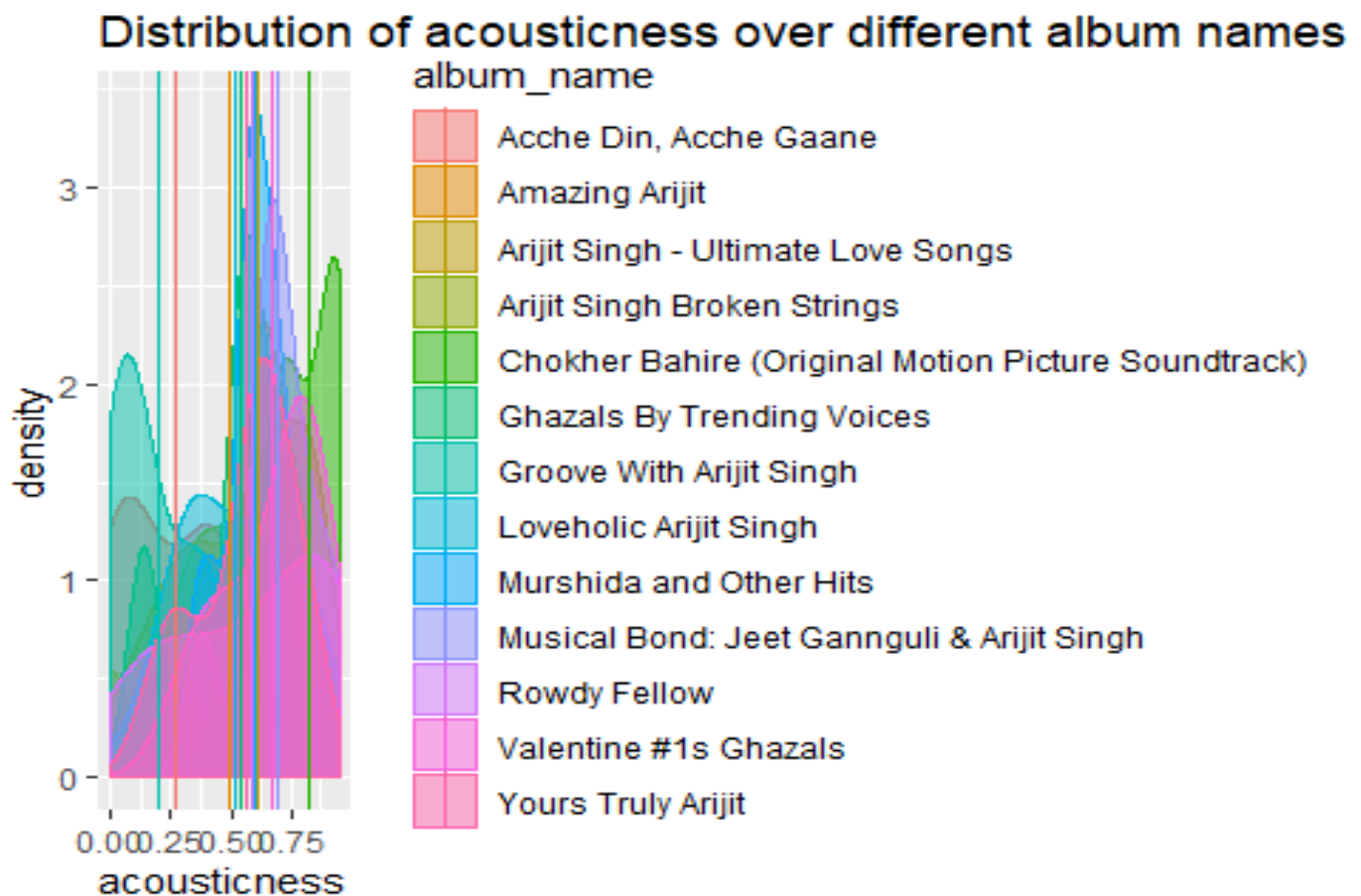
Now doing an ANOVA test to get significance level of getting different group means. Null hypothesis (H_0) states that both the groups have same means.

```
##           Df Sum Sq Mean Sq F value Pr(>F)
## mode_name  1 0.1608 0.160826  6.7231 0.01041 *
```

The results are not significant as expected, since both the group's means could be equal. This is because p-value is close to 0.01 hence not enough evidence to reject null hypothesis. The same information is conveyed by the distribution plot. This can be used to say that energy of his songs cannot be classified using major or minor keys. Energy level of his songs range from low to high for both the keys (i.e. even his minor keys songs can have high energy and major keys songs can have low energy).

- **Pairwise hypothesis tests on acousticness on grouping with album name.**

Distribution of **acousticness** of his songs is studied for each album name. Hence many such different groups are formed.



From the above distribution it can be seen that several groups have similar means (since they are overlapping) while some have different means.

Results of Annova

```
##           Df Sum Sq Mean Sq F value    Pr(>F)
## album_name 12  3.4759  0.289660   5.8366 3.226e-08 ***
```

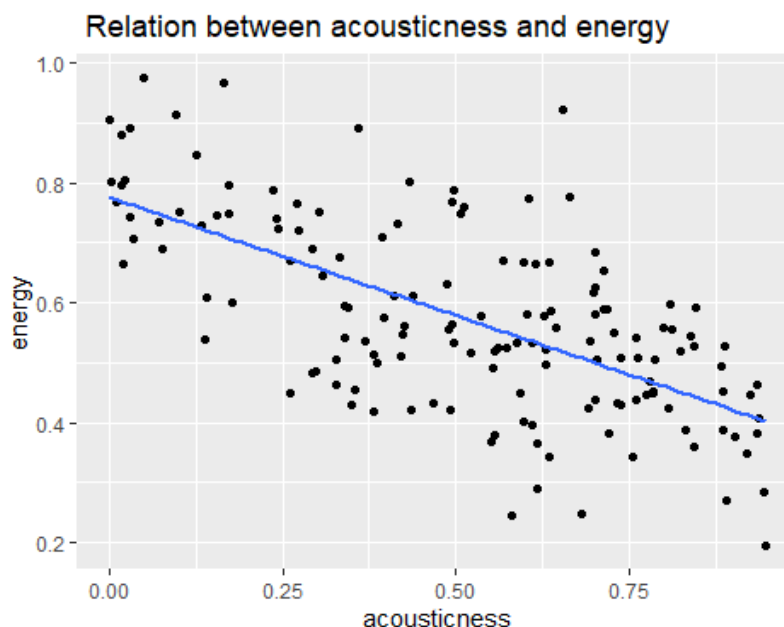
The annova test is significant hence, atleast any one group mean is different. Therefore, digging in deeper to check the pairs of groups which have similar means and which don't.

Upon pairwise **Bonferroni** hypothesis tests we observe that, there are several albums with similar acoustic mean values, but some album groups have different means. Citing an example, albums like 'Grooving with Arjit' and 'Arijit love songs' have very high likely different group means (since p-value almost equal to 0).

This can be explained as 'Grooving with Arjit' has lower acoustic levels as compared to album 'Arijit love songs'.

- **Linear Regression**

Acousticness is chosen as the independent variable and energy of the song as dependent variable. Hence acousticness is predictor and energy as response variable.



The above graph shows a clear linear (decreasing) trend between acousticness and energy of the track.

```
## Multiple R-squared:  0.424, Adjusted R-squared:  0.4204
```

The model tells that around 42% of the data's variance can be explained by the model. Though its not much the model is sufficient to visualize the linear trend between the 2 variables.

As seen earlier the correlation between acousticness and energy is -0.6511881. With increasing acousticness levels the energy of song decreases. This too explains the negative slope of linear fitting curve.

The linear model's summary-

```
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  0.77617    0.02140   36.27  <2e-16 ***
## acousticness -0.39386    0.03652  -10.79  <2e-16 ***
```

It can be mathematically interpreted as:

energy = 0.77617 - 0.39386*acousticness

The y-intercept doesn't have such a physical significance (though its statistically significant with very low p-value) apart from the fact that it acts as a bias for the line to fit the data well.

The slope says energy is dependent on acousticness i.e. with increase in acousticness there is a decrease in energy with a factor of 0.39. Also, the slope and y-intercept are statistically significant since its p-value almost 0 (stating the tru linear model consists such a linear relation between the 2 variables).

The actual slope value does make sense because acousticness and energy are in the same scale ranging from 0 to 1.

Hypothesis Test to check "goodness" of the predictor

A Welch 2 sample test is done to see how the predicted energy levels of track compare to the actual energy levels for those tracks.

Ho = The mean of predicted values and actual values is same.

Ha = The mean of both groups is not same.

```
welch Two Sample t-test
data: fitted(regmod.lm) and df$energy
t = 0, df = 273.29, p-value = 1
alternative hypothesis: true difference in means is not equal to 0
99 percent confidence interval:
 -0.03852419  0.03852419
sample estimates:
mean of x mean of y
 0.5692    0.5692
```

We observe that the mean of both groups is same. Following this is p-value of “1”, hence giving extremely strong evidence of accepting null hypothesis. Thus, means of true and predicted samples is same, showing that ‘acousticness’ is a good predictor for predicting energy level of the songs.

Prediction:

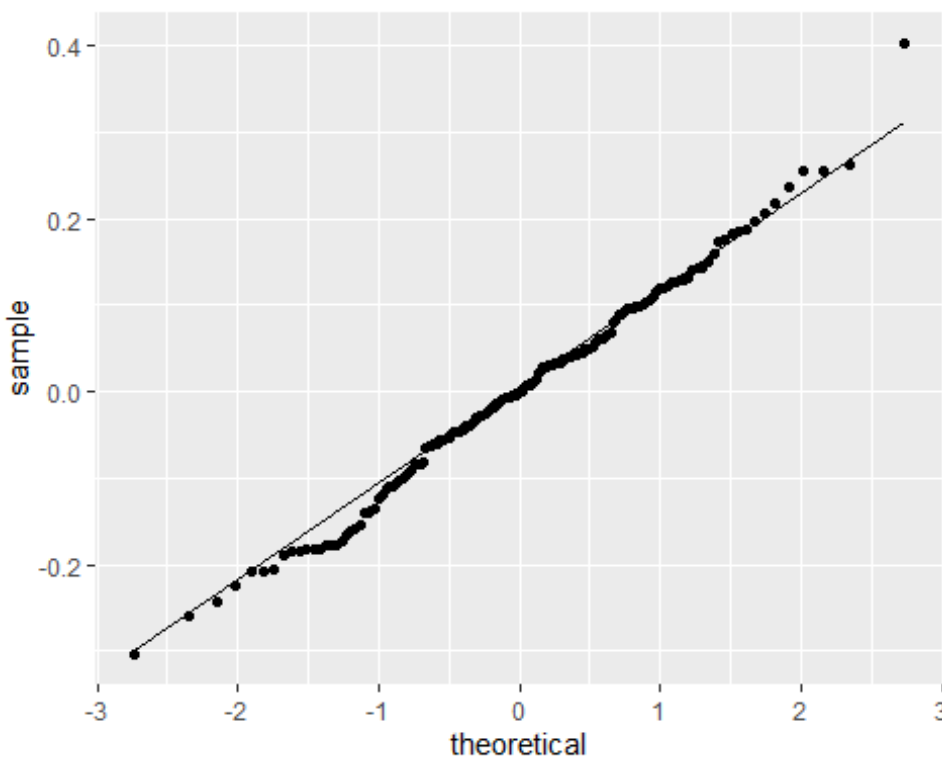
The following is the value of energy levels at the specified percentiles of acousticness levels.

##	20%	40%	60%	80%
##	0.6607659	0.5819149	0.5278771	0.4760449

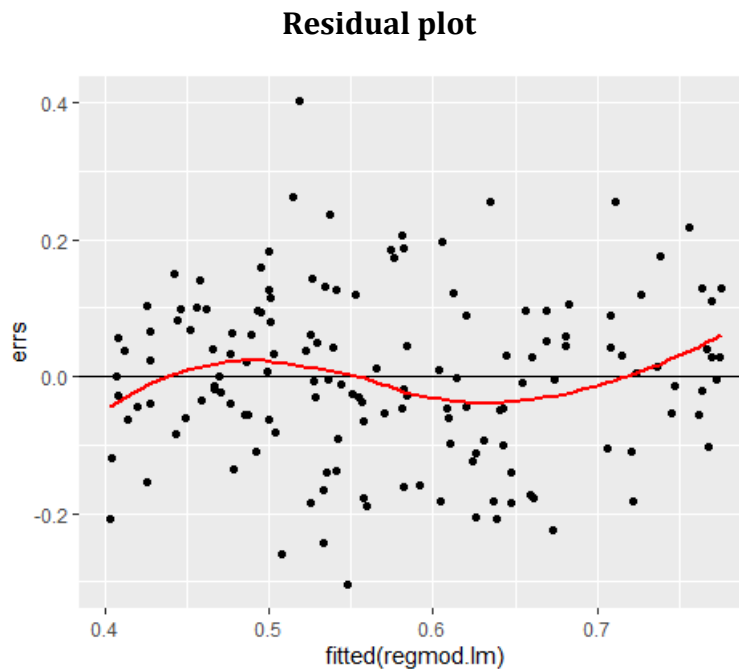
We can see a decreasing trend of energy levels with increasing acousticness levels.

Regression assumptions:

1. Normality



We can see from the above qq-plot that the residuals do form a normal distribution.



2. 0 Mean

The mean of residuals though fluctuates up and down over zero, it can be conveniently concluded that average is almost 0.

3. From the same plot it can be inferred that the variance is almost constant.

Conclusion:

The artist chosen is Arijit Singh and the two most linearly correlated traits are the song's energy levels and its acousticness. Both traits are inversely proportional. Initially it was observed that the track's mode (major or minor) does not much affect the energy of the song since both the modes tend to follow the same distribution. Then, it was observed that different albums of Arijit have different acousticness levels. It was seen that his romantic album has higher acousticness (i.e. lesser energy; might prove his romantic songs are less energizing with less "activity" in the song) and his party album has lesser acousticness (more foot-tapping and energizing song).

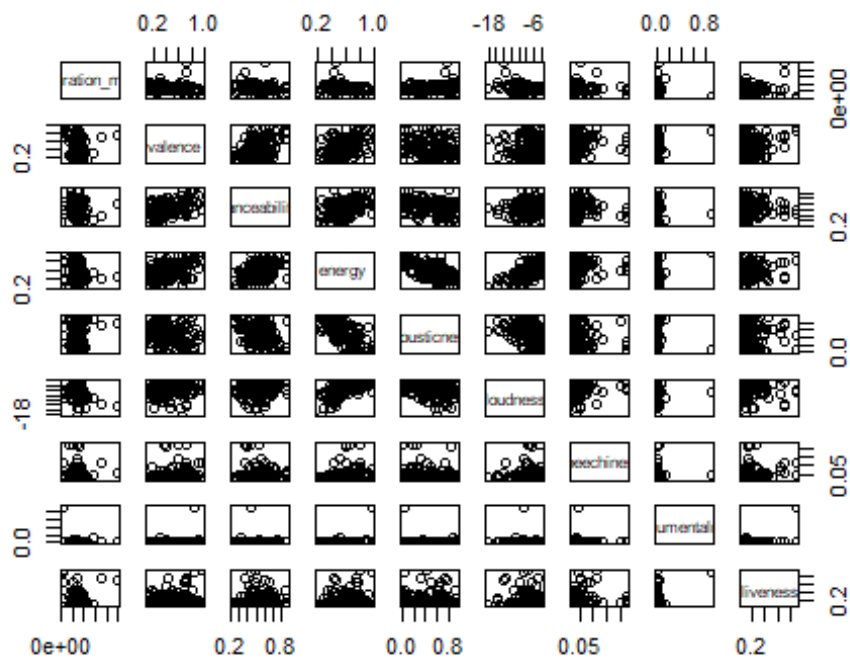
A linear model is fit to predict energy level of a song using acousticness level of that song. The model though doesn't fit the data accurately (most of the variance of data points can't be explained by the model/ a linear curve) but it does successfully show a decreasing trend of energy of the track with increasing acousticness. Hence, most of Arijit's songs that tend to be more acoustic would be less energetic (not a party song but may be a romantic or an emotional song).

APPENDIX
(With R Code)

```
library(spotifyr)
Sys.setenv(SPOTIFY_CLIENT_ID = "b6eb08407bab41b4b3f8f5bb48a0f8f0")
Sys.setenv(SPOTIFY_CLIENT_SECRET = "73efc29aa9b442fdb83019dfc8fc0ed")
access_token <- get_spotify_access_token()
#mmfull <- get_artist_audio_features("modest mouse")

df <- get_artist_audio_features("arjit singh")

pairs(df[c('duration_ms', 'valence', 'danceability', 'energy', 'acousticness',
           'loudness', 'speechiness', 'instrumentalness', 'liveness')])
```



```
cor(df[c('duration_ms', 'valence', 'danceability', 'energy', 'acousticness',
         'loudness', 'speechiness', 'instrumentalness', 'liveness')])
```

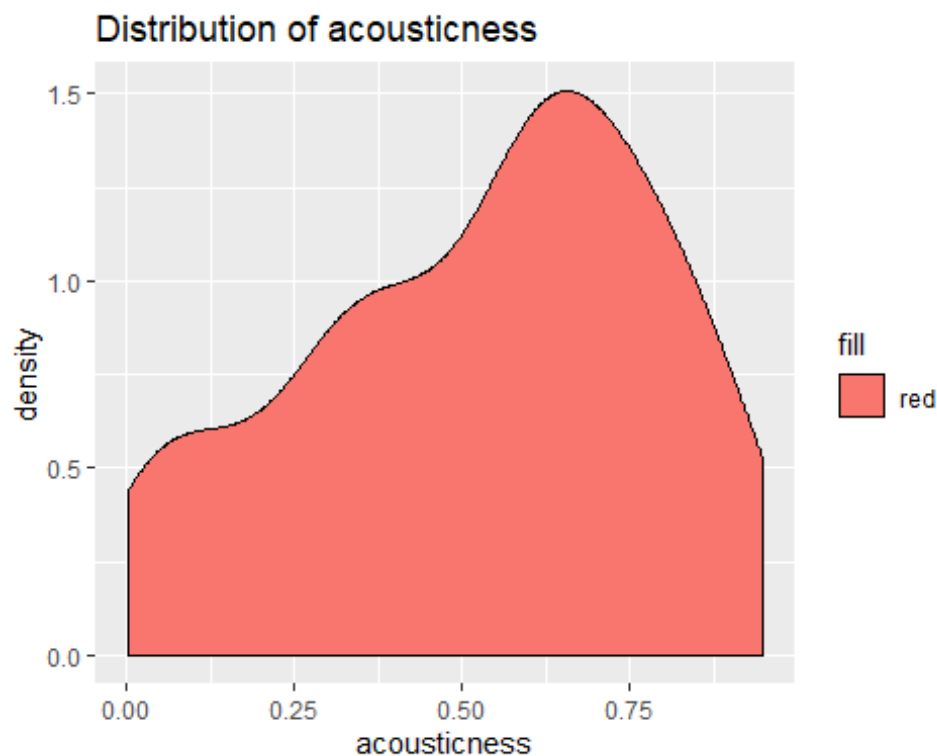
	duration_ms	valence	danceability	energy
duration_ms	1.00000000	-0.16897944	-0.15653746	-0.2071853
valence	-0.16897944	1.00000000	0.40884053	0.4574074
danceability	-0.15653746	0.40884053	1.00000000	0.3476319
energy	-0.20718526	0.45740738	0.34763192	1.00000000
acousticness	0.19638857	-0.27944300	-0.35545567	-0.6511881
loudness	-0.22614206	0.07137788	0.14589959	0.6315025
speechiness	-0.08676018	0.33982239	-0.05737298	0.2936094
instrumentalness	-0.16239510	0.17409380	-0.06840794	0.1952692
liveness	-0.03427193	0.35651840	-0.01692267	0.2390901

```
##
##          acousticness    loudness    speechiness    instrumentalness
## duration_ms      0.1963886 -0.22614206 -0.086760178      -0.162395102
```

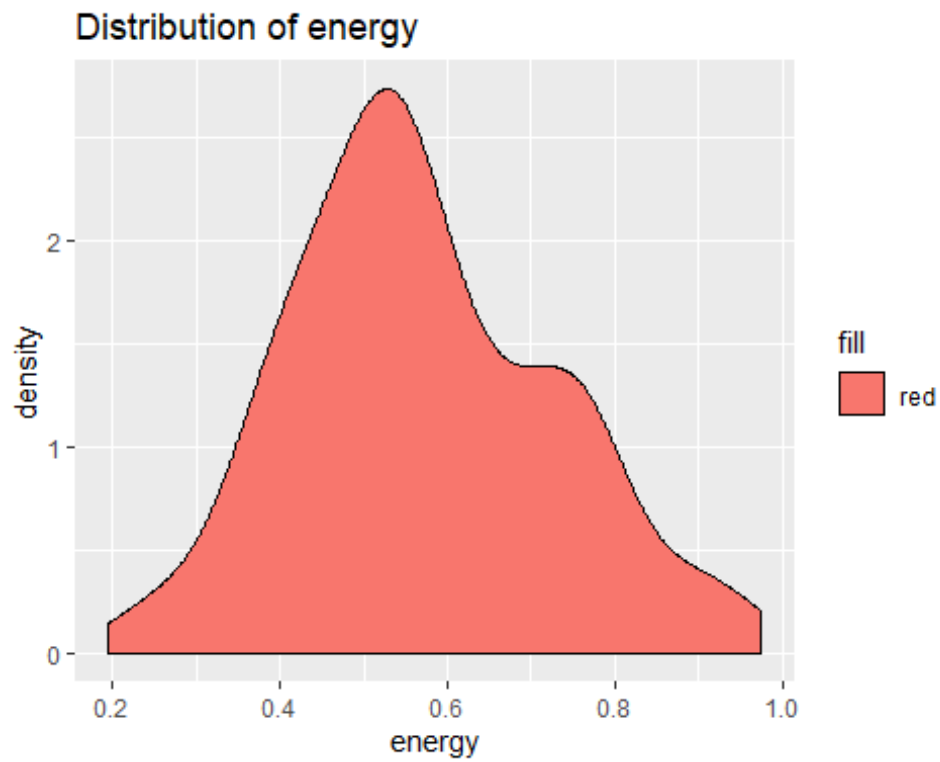
```
## valence          -0.2794430  0.07137788  0.339822390    0.174093803
## danceability     -0.3554557  0.14589959 -0.057372983    -0.068407943
## energy           -0.6511881  0.63150247  0.293609409    0.195269229
## acousticness     1.0000000 -0.34102376 -0.262753114    -0.120519056
## loudness         -0.3410238  1.00000000  0.077353949    -0.057970346
## speechiness      -0.2627531  0.07735395  1.000000000    0.002079317
## instrumentalness -0.1205191 -0.05797035  0.002079317    1.000000000
## liveness         -0.1294266 -0.08826319  0.205140984    0.358387984
##                  liveness
## duration_ms      -0.03427193
## valence           0.35651840
## danceability     -0.01692267
## energy            0.23909008
## acousticness     -0.12942657
## loudness         -0.08826319
## speechiness       0.20514098
## instrumentalness  0.35838798
## liveness          1.00000000
```

Loudness and energy are the most linearly correlated variables

```
ggplot(df, aes(x = acousticness, y = ..density.., fill = 'red')) +
  geom_density() +
  ggtitle("Distribution of acousticness")
```

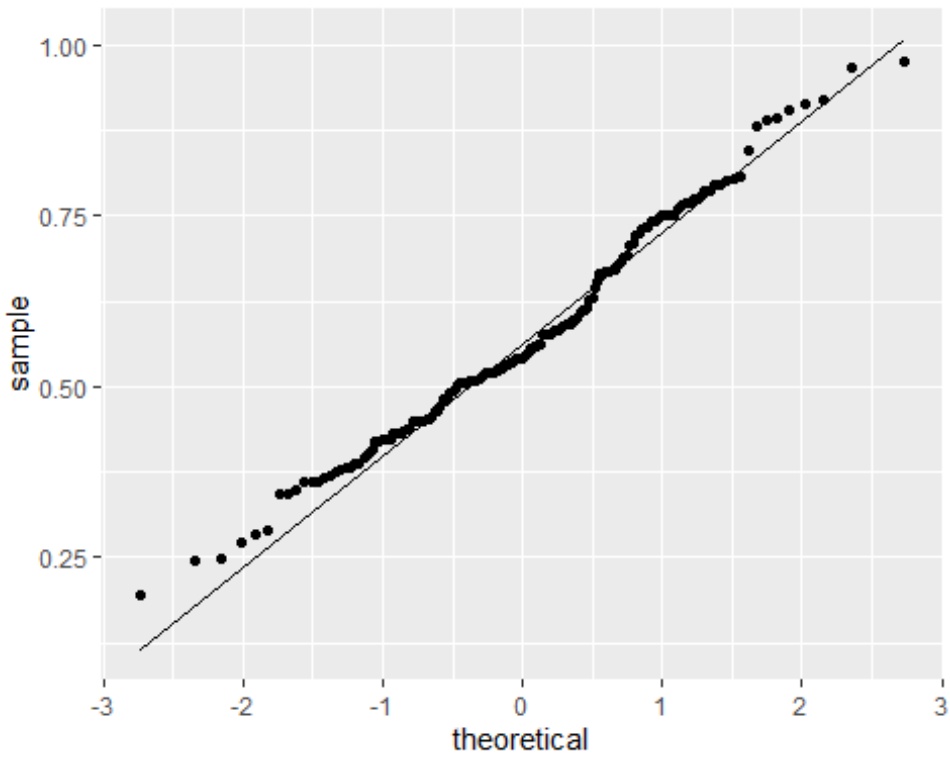


```
ggplot(df, aes(x = energy, y = ..density.., fill = 'red')) +  
geom_density() +  
ggtitle("Distribution of energy")
```



Normality of energy

```
ggplot(df, aes(sample = energy)) + stat_qq() + stat_qq_line()
```

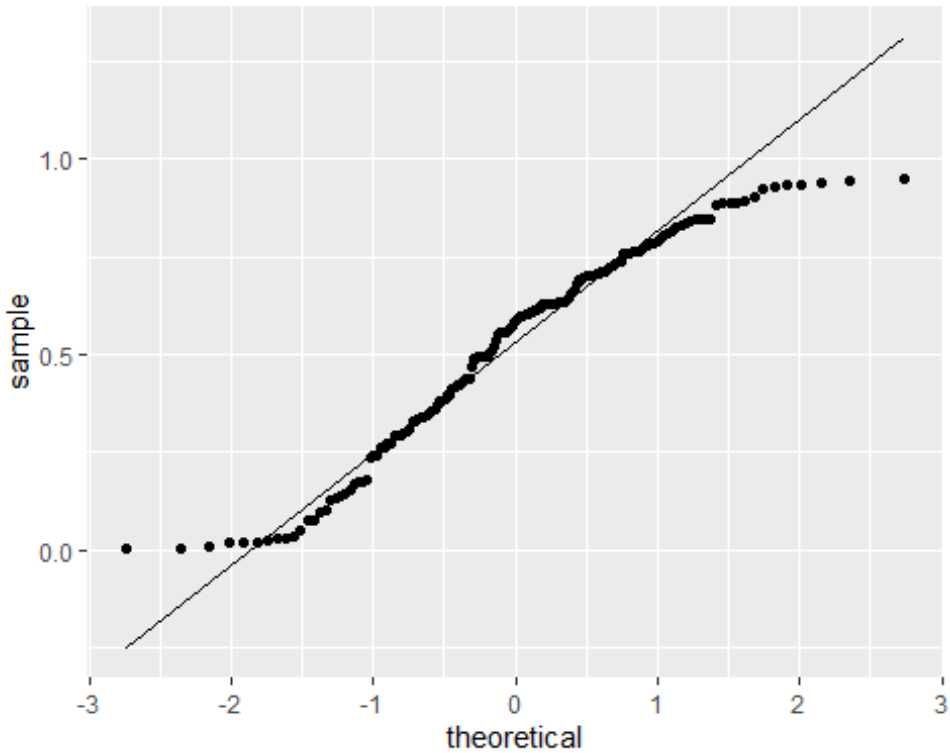


```
ggtitle("Normal QQ plots for energy")
```

```
## $title  
## [1] "Normal QQ plots for energy"  
##  
## attr(,"class")  
## [1] "labels"
```

Normality of acoustiness

```
ggplot(df, aes(sample = acoustiness)) + stat_qq() + stat_qq_line()
```



```
ggtitle("Normal QQ plots for acousticness")

## $title
## [1] "Normal QQ plots for acousticness"
##
## attr(,"class")
## [1] "labels"

CI_energy = t.test(df$energy,conf.level = 0.99)$conf.int
CI_acousticness = t.test(df$acousticness,conf.level = 0.99)$conf.int

print(CI_energy)

## [1] 0.5367533 0.6016467
## attr(,"conf.level")
## [1] 0.99

print(CI_acousticness)

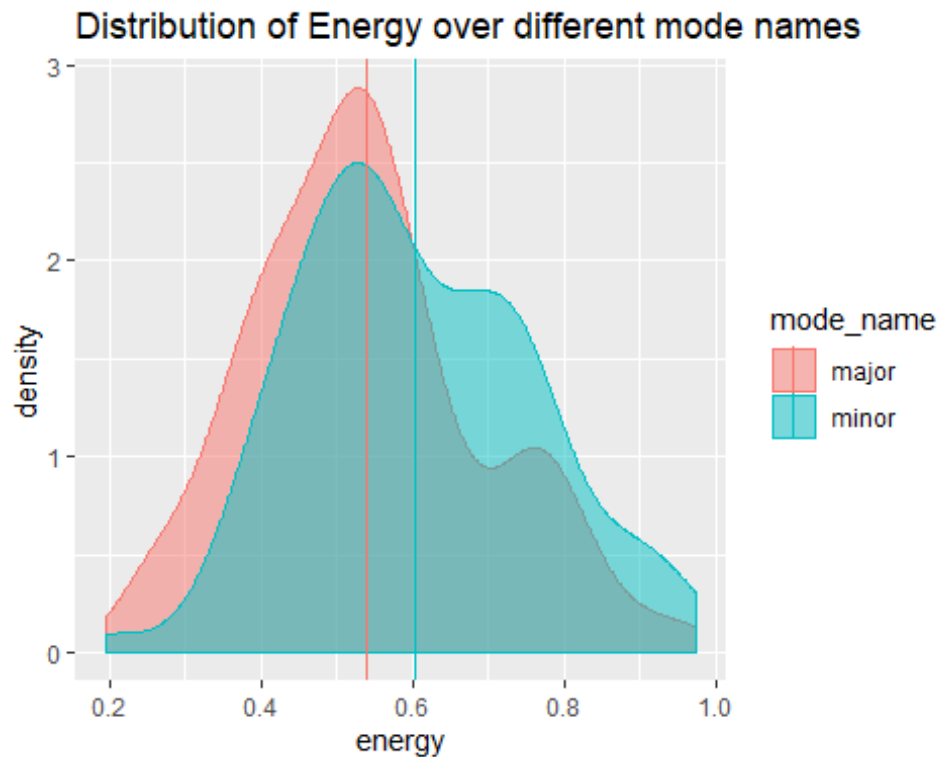
## [1] 0.4718371 0.5791281
## attr(,"conf.level")
## [1] 0.99
```

Q2)

2 sample Hypothesis test on Energy

```
sam.means = aggregate(formula = energy ~ mode_name, data = df, FUN = mean)

ggplot(df, aes(x = energy, y = ..density.., color = mode_name, fill = mode_name)) +
  geom_density(alpha = 0.5) +
  geom_vline(data = sam.means, aes(xintercept = energy, color = mode_name)) +
  ggtitle("Distribution of Energy over different mode names")
```



```
#pairwise.t.test(df$energy, df$mode_name, p.adjust.method = 'bonferroni')
model.lm = lm(energy ~ mode_name, data = df)
anova(model.lm)

## Analysis of Variance Table
##
## Response: energy
##          Df Sum Sq Mean Sq F value Pr(>F)
## mode_name  1  0.1608  0.160826   6.7231 0.01041 *
## Residuals 158  3.7796  0.023921
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

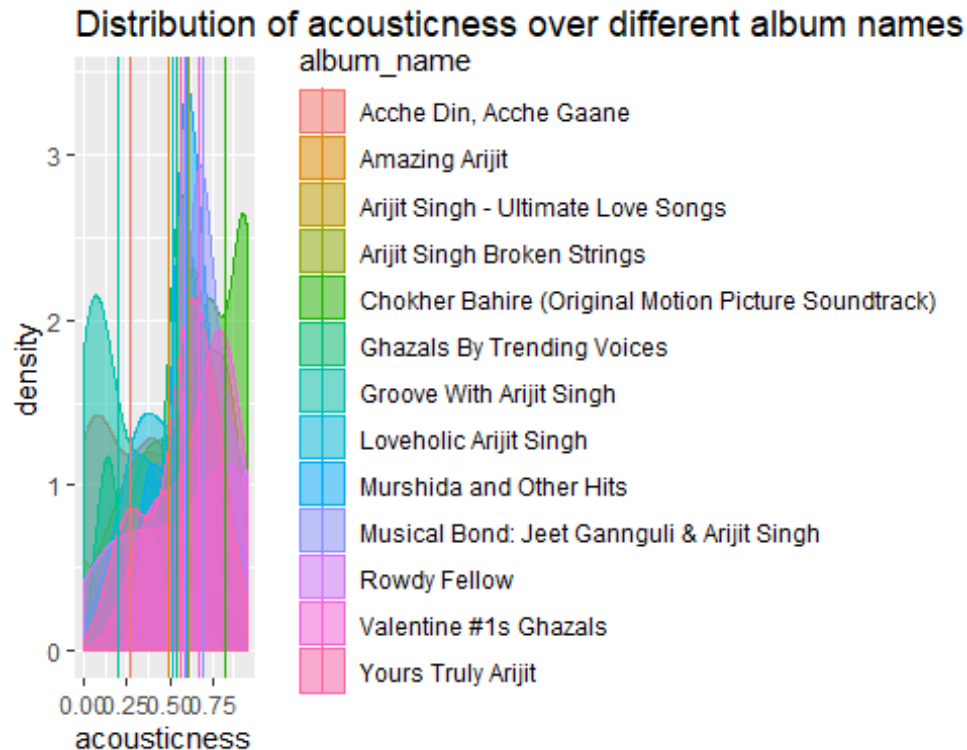
Q3)

pairwise on acousticness on grouping with album name.

```
loud.means = aggregate(formula = acousticness ~ album_name, data = df, FUN = mean)
```



```
ggplot(df, aes(x = acoustiness, y = ..density.., color = album_name, fill = album_name)) +
  geom_density(alpha = 0.5) +
  geom_vline(data = loud.means, aes(xintercept = acoustiness, color = album_name)) +
  ggtitle("Distribution of acoustiness over different album names")
```



```
model.lm = lm(acoustiness ~ album_name, data = df)
anova(model.lm)

## Analysis of Variance Table
##
## Response: acoustiness
##           Df Sum Sq Mean Sq F value    Pr(>F)
## album_name 12  3.4759  0.289660   5.8366 3.226e-08 ***
## Residuals 147  7.2954  0.049628
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

The anova test is significant hence, atleast any one group mean is different. Therefore, digging in deeper to check the pairs of groups which have similar means and which don't.

```
pairwise.t.test(df$acoustiness, df$album_name, p.adjust.method = 'bonferroni')
')
```

```

##
## Pairwise comparisons using t tests with pooled SD
##
## data: df$acousticness and df$album_name
##
##                                     Acche Din, Acche Gaane
## Amazing Arijit                      1.00000
## Arijit Singh - Ultimate Love Songs  0.04074
## Arijit Singh Broken Strings          0.08336
## Chokher Bahire (Original Motion Picture Soundtrack) 0.00013
## Ghazals By Trending Voices          0.45394
## Groove With Arijit Singh            1.00000
## Loveholic Arijit Singh              0.21547
## Murshida and Other Hits             0.06879
## Musical Bond: Jeet Gannguli & Arijit Singh 0.00365
## Rowdy Fellow                       0.06525
## Valentine #1s Ghazals               0.00353
## Yours Truly Arijit                 0.03559
##                                     Amazing Arijit
## Amazing Arijit                      -
## Arijit Singh - Ultimate Love Songs  1.00000
## Arijit Singh Broken Strings          1.00000
## Chokher Bahire (Original Motion Picture Soundtrack) 0.38928
## Ghazals By Trending Voices          1.00000
## Groove With Arijit Singh            0.18695
## Loveholic Arijit Singh              1.00000
## Murshida and Other Hits             1.00000
## Musical Bond: Jeet Gannguli & Arijit Singh 1.00000
## Rowdy Fellow                       1.00000
## Valentine #1s Ghazals               1.00000
## Yours Truly Arijit                 1.00000
##                                     Arijit Singh - Ultimat
e Love Songs
## Amazing Arijit                      -
## Arijit Singh - Ultimate Love Songs  -
## Arijit Singh Broken Strings          1.00000
## Chokher Bahire (Original Motion Picture Soundtrack) 1.00000
## Ghazals By Trending Voices          1.00000
## Groove With Arijit Singh            0.00133
## Loveholic Arijit Singh              1.00000
## Murshida and Other Hits             1.00000
## Musical Bond: Jeet Gannguli & Arijit Singh 1.00000
## Rowdy Fellow                       1.00000
## Valentine #1s Ghazals               1.00000
## Yours Truly Arijit                 1.00000
##                                     Arijit Singh Broken St
rings
## Amazing Arijit                      -
## Arijit Singh - Ultimate Love Songs  -
## Arijit Singh Broken Strings          -

```

## Chokher Bahire (Original Motion Picture Soundtrack)	0.52006
## Ghazals By Trending Voices	1.00000
## Groove With Arijit Singh	0.00147
## Loveholic Arijit Singh	1.00000
## Murshida and Other Hits	1.00000
## Musical Bond: Jeet Gannguli & Arijit Singh	1.00000
## Rowdy Fellow	1.00000
## Valentine #1s Ghazals	1.00000
## Yours Truly Arijit	1.00000
## Chokher Bahire (Original Motion Picture Soundtrack)	-
## Amazing Arijit	-
## Arijit Singh - Ultimate Love Songs	-
## Arijit Singh Broken Strings	-
## Chokher Bahire (Original Motion Picture Soundtrack)	-
## Ghazals By Trending Voices	0.99714
## Groove With Arijit Singh	3.3e-06
## Loveholic Arijit Singh	0.26119
## Murshida and Other Hits	1.00000
## Musical Bond: Jeet Gannguli & Arijit Singh	1.00000
## Rowdy Fellow	1.00000
## Valentine #1s Ghazals	1.00000
## Yours Truly Arijit	0.90097
## Ghazals By Trending Voices	-
## Amazing Arijit	-
## Arijit Singh - Ultimate Love Songs	-
## Arijit Singh Broken Strings	-
## Chokher Bahire (Original Motion Picture Soundtrack)	-
## Ghazals By Trending Voices	-
## Groove With Arijit Singh	0.02611
## Loveholic Arijit Singh	1.00000
## Murshida and Other Hits	1.00000
## Musical Bond: Jeet Gannguli & Arijit Singh	1.00000
## Rowdy Fellow	1.00000
## Valentine #1s Ghazals	1.00000
## Yours Truly Arijit	1.00000
## Groove With Arijit Singh	-
## Amazing Arijit	-
## Arijit Singh - Ultimate Love Songs	-
## Arijit Singh Broken Strings	-
## Chokher Bahire (Original Motion Picture Soundtrack)	-
## Ghazals By Trending Voices	-
## Groove With Arijit Singh	-
## Loveholic Arijit Singh	0.00500
## Murshida and Other Hits	0.00251
## Musical Bond: Jeet Gannguli & Arijit Singh	0.00010
## Rowdy Fellow	0.00211
## Valentine #1s Ghazals	7.1e-05

## Yours Truly Arijit	0.00050
##	Loveholic Arijit Singh
## Amazing Arijit	-
## Arijit Singh - Ultimate Love Songs	-
## Arijit Singh Broken Strings	-
## Chokher Bahire (Original Motion Picture Soundtrack)	-
## Ghazals By Trending Voices	-
## Groove With Arijit Singh	-
## Loveholic Arijit Singh	-
## Murshida and Other Hits	1.00000
## Musical Bond: Jeet Gannguli & Arijit Singh	1.00000
## Rowdy Fellow	1.00000
## Valentine #1s Ghazals	1.00000
## Yours Truly Arijit	1.00000
##	Murshida and Other Hit
S	
## Amazing Arijit	-
## Arijit Singh - Ultimate Love Songs	-
## Arijit Singh Broken Strings	-
## Chokher Bahire (Original Motion Picture Soundtrack)	-
## Ghazals By Trending Voices	-
## Groove With Arijit Singh	-
## Loveholic Arijit Singh	-
## Murshida and Other Hits	-
## Musical Bond: Jeet Gannguli & Arijit Singh	1.00000
## Rowdy Fellow	1.00000
## Valentine #1s Ghazals	1.00000
## Yours Truly Arijit	1.00000
##	Musical Bond: Jeet Gan
nguli & Arijit Singh	
## Amazing Arijit	-
## Arijit Singh - Ultimate Love Songs	-
## Arijit Singh Broken Strings	-
## Chokher Bahire (Original Motion Picture Soundtrack)	-
## Ghazals By Trending Voices	-
## Groove With Arijit Singh	-
## Loveholic Arijit Singh	-
## Murshida and Other Hits	-
## Musical Bond: Jeet Gannguli & Arijit Singh	-
## Rowdy Fellow	1.00000
## Valentine #1s Ghazals	1.00000
## Yours Truly Arijit	1.00000
##	Rowdy Fellow
## Amazing Arijit	-
## Arijit Singh - Ultimate Love Songs	-
## Arijit Singh Broken Strings	-
## Chokher Bahire (Original Motion Picture Soundtrack)	-
## Ghazals By Trending Voices	-
## Groove With Arijit Singh	-
## Loveholic Arijit Singh	-

```

## Murshida and Other Hits -
## Musical Bond: Jeet Gannguli & Arijit Singh -
## Rowdy Fellow -
## Valentine #1s Ghazals 1.00000
## Yours Truly Arijit 1.00000
## Valentine #1s Ghazals
## Amazing Arijit -
## Arijit Singh - Ultimate Love Songs -
## Arijit Singh Broken Strings -
## Chokher Bahire (Original Motion Picture Soundtrack) -
## Ghazals By Trending Voices -
## Groove With Arijit Singh -
## Loveholic Arijit Singh -
## Murshida and Other Hits -
## Musical Bond: Jeet Gannguli & Arijit Singh -
## Rowdy Fellow -
## Valentine #1s Ghazals -
## Yours Truly Arijit 1.00000
##
## P value adjustment method: bonferroni

```

There are several albums with similar loudness mean values but some album groups have different means.

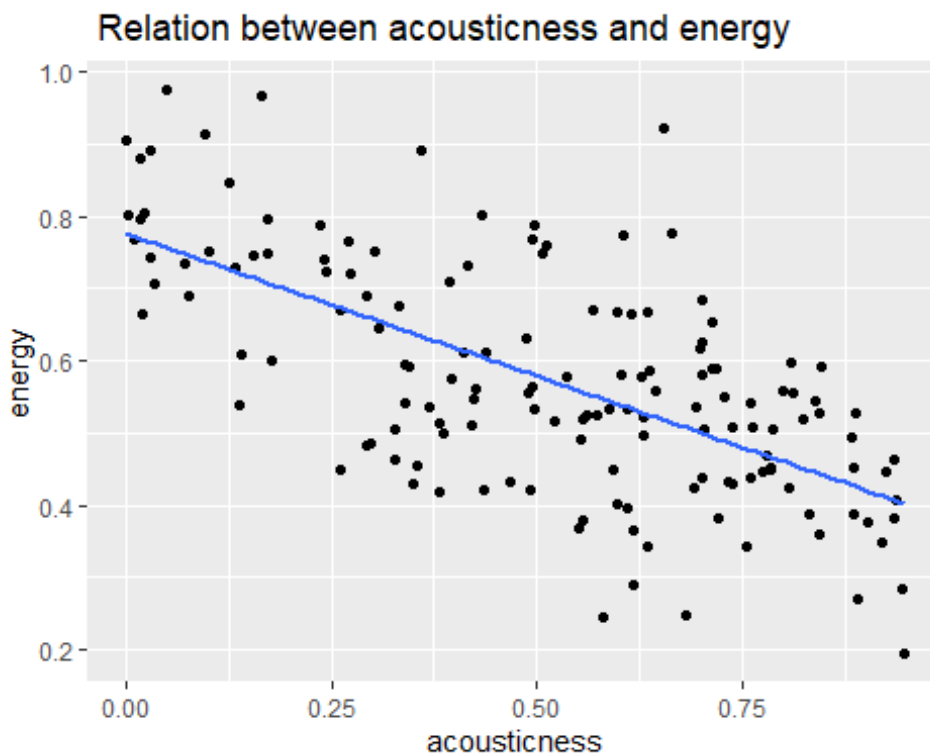
Q4) Linear Regression

acousticness is chosen as the independent variable and energy of the song as dependent variable. Hence acousticness is predictor and energy as response variable.

```

ggplot(df, aes(x = acousticness, y = energy)) +
  geom_point() +
  geom_smooth(method = 'lm', se = FALSE) +
  ggtitle(" Relation between acousticness and energy")

```



```
print(cor(df$acousticness, df$energy))
```

```
## [1] -0.6511881
```

Building a linear regression model as above with acousticness as the independent variable and energy of the song as dependent variable.

```
regmod.lm = lm(energy ~ acousticness, df)
```

```
summary(regmod.lm)
```

```
##
## Call:
## lm(formula = energy ~ acousticness, data = df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.30373 -0.06973 -0.00112  0.08112  0.40281
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   0.77617    0.02140   36.27  <2e-16 ***
## acousticness -0.39386    0.03652  -10.79  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.1198 on 158 degrees of freedom
```

```
## Multiple R-squared:  0.424, Adjusted R-squared:  0.4204
## F-statistic: 116.3 on 1 and 158 DF,  p-value: < 2.2e-16
```

c) H_0 = The mean of predicted values and actual values is same.

H_a = The mean of both groups is not same.

```
```{r}
```

```
t.test(fitted(regmod.lm), df$energy, conf.level = 0.99)
```

```
```
```

```
welch Two Sample t-test
```

```
data: fitted(regmod.lm) and df$energy
t = 0, df = 273.29, p-value = 1
alternative hypothesis: true difference in means is not equal to 0
99 percent confidence interval:
 -0.03852419  0.03852419
sample estimates:
mean of x mean of y
 0.5692    0.5692
```

```
## generating test points as specific percentiles of train data
test_points = quantile(df$acousticness, probs = c(0.2, 0.4, 0.6, 0.8))
```

```
# predicting energy values using those loudness data points
predict(regmod.lm, data.frame(acousticness = test_points))
```

```
##          20%          40%          60%          80%
## 0.6607659 0.5819149 0.5278771 0.4760449
```

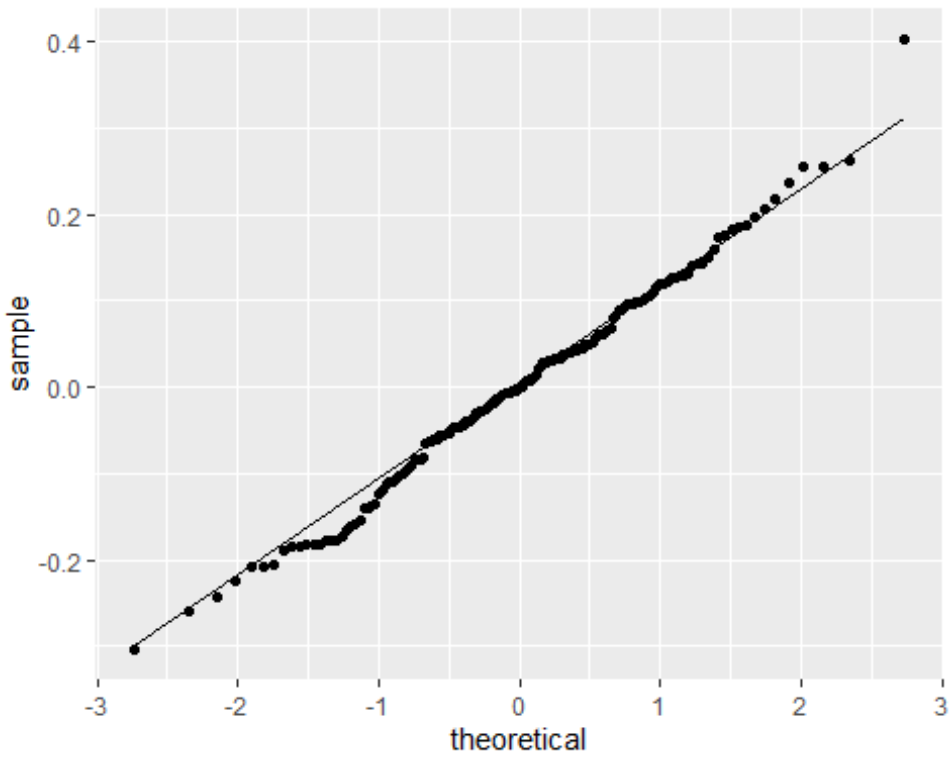
e) Regression assumptions

```
# getting residuals
```

```
errs= residuals(regmod.lm)
```

```
## normality of residulas ##
```

```
ggplot(df, aes(sample = errs)) + stat_qq() + stat_qq_line()
```



```
## residuals plot ##
```

```
ggplot(df, aes(fitted(regmod.lm), errs)) + geom_point() + geom_hline(yintercept= 0) +  
  geom_smooth(col = 'red', se = FALSE)
```

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