

Intro	Context	Artist Panorama	Artist Trends	Datasets and Processes	Models_01	Models_02	Models_03	Conclusions
-------	---------	-----------------	---------------	------------------------	-----------	-----------	-----------	-------------



Data Analytics Bootcamp Project 04

Team 01:
Jaime Bravo
Roberto Rodas..

Intro	Context	Artist Panorama	Artist Trends	Datasets and Processes	Models_01	Models_02	Models_03	Conclusions
-------	---------	-----------------	---------------	------------------------	-----------	-----------	-----------	-------------



Context

Re-exploring the topic of our first project, we decided to apply more sophisticated tools to develop models that could predict success for a musical project today, considering the vast amounts of data generated by musical streaming and traditional data-collection means, in order to create valuable information for industry executives, musicians, marketing teams and ..

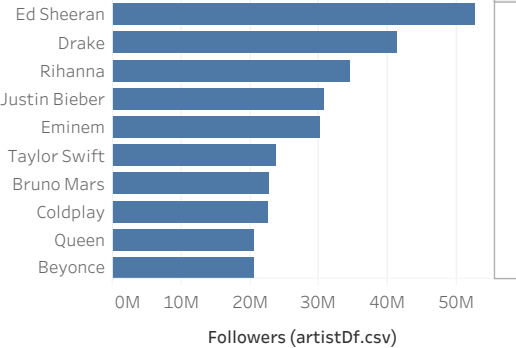
Research Questions

- How are stream numbers related to other success metrics, such as certifications and chart positioning.
- How do song attributes relate to musical success?
- Which labels produce more successful releases?
- Does featuring with another artist raise the success probabilities?

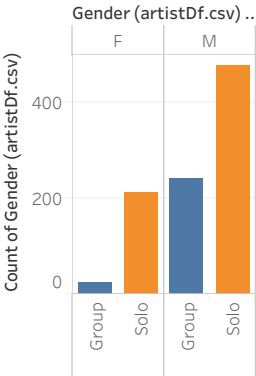
Artist Panorama

We considered it was important to display the general panorama of the data that we worked with, and the element that was at the center of it a..

Artist (artistDf.csv)



The data that we worked with comprised all entries into the Billboard Charts from the year 1999 to 2019. In this first table we can see artists ranked by their number o..

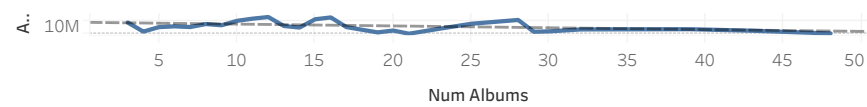


And in this table we can see how the genders, as well as type of aggrupation are distributed. We can see more prevalence of male artists, and a greater tendency ..

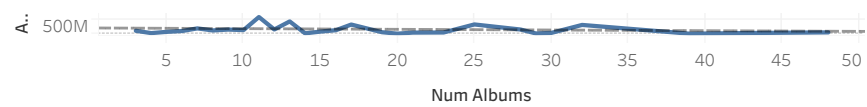
Artist Trends

We also identified some trends regarding the number of albums and the popularity measured in average streams and followers, as well as the p..

Albums-Followers



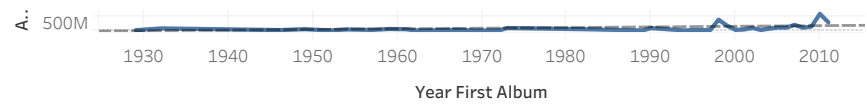
Albums-Streams



Year-Followers



Year-Streams



This simple analysis showed us that there were relations worth exploring in relation to artist popularity, not only relating to the albums' release, but to the music itself as well.

Intro	Context	Artist Panorama	Artist Trends	Datasets and Processes	Models_01	Models_02	Models_03	Conclusions
-------	---------	-----------------	---------------	------------------------	-----------	-----------	-----------	-------------

Datasets and Processes

Considering these findings, we continued with the processing of data in order to create machine learning models that could help predict the pop..

The CSVs that compose our oriignal data, all coming from Billboard charts, and the Spotify API were:

- artistDf
- billboardHot100_1999-2019
- grammyAlbums_1999-2019
- grammySongs_1999-2019
- riaaAlbumCerts_1999-2019..

Which we transformed into:

- artistWeek, by merging artistDf and spotifyWeeklyTop200Streams
- attributesBillboard, by merging songAttributes_1999-2019 and billboard..



To apply the following processes:

- Brainstorming and developing research questions
- Data cleaning and preprocessing
- Model testing
- SQLite integration
- Gathering results..

Using tools such as:

- Python
- Pandas
- Numpy
- SQLite..

Intro	Context	Artist Panorama	Artist Trends	Datasets and Processes	Models_01	Models_02	Models_03	Conclusions
-------	---------	-----------------	---------------	------------------------	-----------	-----------	-----------	-------------

Models

Here are the results from the models we used

01.- Neural Network: Attributes Sample

Model: "sequential"

Layer (type)	Output Shape	Param #
Dense (Dense)	(None, 20)	200
Dense_1 (Dense)	(None, 20)	420
Dense_2 (Dense)	(None, 1)	21

total params: 641 (2.56 KB)
 trainable params: 641 (2.56 KB)
 non-trainable params: 0 (0.00 Byte)

Epochs: 30

Function: Relu/Sigmoid

This model evaluated popularity exclusively agai..

```
# Evaluate the performance of model using the loss and predictive accuracy of the model on the test dataset.
model_loss, model_accuracy = m_model.evaluate(X_test_scaled,y_test,verbose=2)
print(f"Loss: {model_loss}, Accuracy: {model_accuracy}")

✓ file

364/364 - 3s - loss: 0.5272 - accuracy: 0.7614 - 580ms/epoch - 2ms/step
loss: 0.5271784661280378, Accuracy: 0.7614231117786299
```

01.- Neural Network: AttributesBillboard

Model: "sequential_1"

Layer (type)	Output Shape	Param #
Dense_1 (Dense)	(None, 40)	160
Dense_2 (Dense)	(None, 40)	1640
Dense_3 (Dense)	(None, 40)	1640
Dense_4 (Dense)	(None, 1)	41

total params: 3881 (25.96 KB)
 trainable params: 3881 (25.96 KB)
 non-trainable params: 0 (0.00 Byte)

Epochs: 50

Function: Relu/Sigmoid

This model evaluated popularity against song att..

```
# Evaluate the model using the test data
model_loss, model_accuracy = m_model.evaluate(X_test_scaled,y_test,verbose=2)
print(f"Loss: {model_loss}, Accuracy: {model_accuracy}")

✓ 0.2s

36/36 - 8s - loss: 0.6803 - accuracy: 0.6803 - 120ms/epoch - 5ms/step
loss: 0.6837423140766663, Accuracy: 0.68255388884892
```

The metrics that we used to evaluate success in all of our models were either popularity or number of followers, and for binary classification purposes, we defined a release as "is_popular" when it had a value equal or greater than the value at the 75% quartile. Also, the epoch and activation function selection came after iterating different options, and sel..

Models

Here are the results from the models we used

03.- Neural Network: ArtistWeek

Layer (type)	Output Shape	Param #
dense_18 (Dense)	(None, 96)	276
dense_19 (Dense)	(None, 96)	9168
dense_20 (Dense)	(None, 96)	9168
dense_21 (Dense)	(None, 1)	96

Total params: 1941 (3.46 M)
Trainable params: 1941 (3.46 M)
Non-trainable params: 0 (0.00 M)

Epochs: 50
Function: Relu/Sigmoid

This model evaluated popularity by the amount o..

```
# Evaluate the model using the test data
model_loss, model_accuracy = nn.evaluate(X_test_scaled, y_test, verbose=2)
print('Loss: {model_loss}, Accuracy: {model_accuracy}')
✓ Run

2/2 - loss: 0.4504 - accuracy: 0.8445 - 50ms/epoch - 75ms/step
loss: 0.450426110308027, Accuracy: 0.8444668717716675
```

04.- Neural Network: Attributes Full

Layer (type)	Output Shape	Param #
dense_18 (Dense)	(None, 96)	276
dense_19 (Dense)	(None, 96)	9168
dense_20 (Dense)	(None, 96)	9168
dense_21 (Dense)	(None, 1)	96

Total params: 1941 (3.46 M)
Trainable params: 1941 (3.46 M)
Non-trainable params: 0 (0.00 M)

Epochs: 30
Function: Relu/Sigmoid

Similar to the first model, this one evaluates pop..

```
# Evaluate the model using the test data
model_loss, model_accuracy = nn.evaluate(X_test_scaled, y_test, verbose=2)
print('Loss: {model_loss}, Accuracy: {model_accuracy}')
✓ Run

1211/1211 - loss: 0.5700 - accuracy: 0.7824 - 2s/epoch - 28ms/step
loss: 0.570736647888227, Accuracy: 0.78240621111097
```

The metrics that we used to evaluate success in all of our models were either popularity or number of followers, and for binary classification purposes, we defined a release as "is_popular" when it had a value equal or greater than the value at the 75% quartile. Also, the epoch and activation function selection came after iterating different options, and sel..

Intro	Context	Artist Panorama	Artist Trends	Datasets and Processes	Models_01	Models_02	Models_03	Conclusions
-------	---------	-----------------	---------------	------------------------	-----------	-----------	-----------	-------------

Models

Here are the results from the models we used

05.- SVC: AttributesBillboard

```

> calculate the classification report
predictions = model.predict(X_test)
print(classification_report(y_test, predictions,
                           target_names=target.popularity))
✓ 0/0

```

	precision	recall	f1 score	support
popular	0.76	1.00	0.87	993
not popular	0.00	0.00	0.00	283
accuracy	0.76			
macro avg	0.38	0.50	0.43	1186
weighted avg	0.76	0.76	0.66	1186

This model evaluated popularity against song attributes, as well as Billboard Charts information, but used an SVC Model instead of Neural Networks.

```

print("training data score: (model.score(X_train, y_train))")
print("testing data score: (model.score(X_test, y_test))")
✓ 0/0

```

Training Data Score:	0.7625698124822346
Testing Data Score:	0.7629015745393634

06.- Random Forest: Grammy/Label

```

Confusion Matrix

```

	Predicted 0	Predicted 1
Actual 0	75	51
Actual 1	42	107

```

Accuracy Score: 0.65580734067297

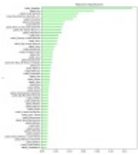
```

```

Classification Report

```

	precision	recall	f1 score	support
label	0.63	0.46	0.54	126
label	0.68	0.75	0.72	159
accuracy	0.63			
macro avg	0.65	0.61	0.63	285
weighted avg	0.63	0.61	0.60	285



This model, the most different one, didn't evaluate popularity or number of followers, instead it used a Random Forest Model to evaluate the likelihood of getting a grammy in relation to music publishing labels.

The metrics that we used to evaluate success in all of our models were either popularity or number of followers, and for binary classification purposes, we defined a release as "is_popular" when it had a value equal or greater than the value at the 75% quartile. Also, the epoch and activation function selection came after iterating different options, and sel..

Intro	Context	Artist Panorama	Artist Trends	Datasets and Processes	Models_01	Models_02	Models_03	Conclusions
-------	---------	-----------------	---------------	------------------------	-----------	-----------	-----------	-------------

Conclusions

- Neural Networks seem to have the highest accuracy scores for these types of evaluations
- Spending considerable effort on cleaning up the data, and performing feature engineering tends to have a much greater impact on model accuracy
- We found that increasing epochs over 50, having more than 3 layers or having more than 30 neurons per layer didn't seem to have a more positive impact on accuracy on any of the tests.
- For our data, it seems that there is a reasonable correlation between musical attributes and success of a musical release, as well as the specifics of the artist, their career history, gender and aggrupation type.
- We have to be very critical in evaluating which features contribute to our target, such as popularity, and which features are consequences of said popularity, and might be impacted by

