

# Clustering of Songs with Amazon Web Services



Certification project  
for AIDA 2020

Falk Lutz  
Julian Godley  
Simon Harston

# Table of Contents

<b>Table of Contents</b>	<b>2</b>
<b>List of Abbreviations</b>	<b>2</b>
<b>Assignment</b>	<b>3</b>
<b>Prerequisites</b>	<b>3</b>
Technical infrastructure	3
Data availability	3
<b>Process workflow / pipeline</b>	<b>4</b>
Data preprocessing	4
Exploratory Data Analysis (EDA)	4
Clustering Algorithms	7
Distributed Cluster-Computing Frameworks	7
AWS Cloud	7
AWS Sagemaker Standalone Notebook	7
AWS Sagemaker Cluster	9
Apache Spark on Databricks	9
<b>Learnings / Insights</b>	<b>10</b>
Clustering and Recommendation	10
Performance	10
<b>Potential next steps</b>	<b>11</b>
<b>Distribution of tasks</b>	<b>11</b>
<b>List of links</b>	<b>11</b>
Kaggle Challenge site	11
Spotify	11
Clustering Algorithms	12
Kmeans	12
DBScan	12

## List of Abbreviations

AWS	Amazon Web Services
CSV	“Comma separated values”, a format where data columns are separated by comma
DBSCAN	“Density-based spatial clustering of applications with noise”, a clustering method
EDA	Exploratory Data Analysis
PCA	Principal Components Analysis, a method of reducing feature dimensions
RDS	Relational Database Service (part of Amazon Web Services)
S3	Simple Storage Service (part Amazon of Amazon Web Services)

# Assignment

The objective of this project is to group a set of songs allocated in Spotify that have some similarities among them. This dataset is composed of three files retrieved from Spotify with audio features for over 20k songs. Every file contains 18 attributes among which are strings and integer attributes such as artist name, track id, danceability, energy, etc.

**The final goal is to find out new songs that can be part of your list of favourite songs. To accomplish this task, different ways of using a clustering approach are going to be compared. On the one hand, through a clustering approach following a distributed cluster-computing framework, and on the other hand, through the use of Amazon RDS.**

To do that, some tasks need to be accomplished:

1. Understand the content that is available in the dataset and clean the dataset if it is necessary.
2. Clustering following a distributed cluster-computing framework.
  - a. Choose the best clustering algorithm taking into account the attributes that the dataset offers. In order to decide the best approach, it could be useful to have an overview of the different perspectives explaining the advantages and drawbacks of your decision.
  - b. Evaluate your model statistically.
3. Clustering in the cloud using AWS.
  - a. Create a database using Amazon RDS.
  - b. Create tables to load the data of the dataset.
  - c. Export the data into S3 and apply a cluster analysis in AWS using a different clustering algorithm.
  - d. Analyze the outcomes.
  - e. (Optional task) Put the final results into RDS for future access.
4. Use the different packages of visualization explained during the course to visualize clusters based on your findings.
5. Compare the findings of both methods using metrics, the visualizations, etc.

## Prerequisites

### Technical infrastructure

The project team agreed to utilize the same technical infrastructure as had been utilized for a previous project. This included:

- Data Storage on an AWS S3 Bucket (<s3://flutz-bucket/spotify/>)
- Database on AWS RDS ([fjs-project.cpfawmiirogo.eu-central-1.rds.amazonaws.com](https://fjs-project.cpfawmiirogo.eu-central-1.rds.amazonaws.com))
- Coding platforms
  - AWS Sagemaker Notebook Instance ([https://flutz.notebook.eu-central-1.sagemaker.aws/tree/CoS\\_AWS](https://flutz.notebook.eu-central-1.sagemaker.aws/tree/CoS_AWS))
  - Databricks workspace (<https://community.cloud.databricks.com/?o=6620618385044013>)
  - Repository: GitHub ([https://github.com/jgo-DA/CoS\\_AWS](https://github.com/jgo-DA/CoS_AWS))

### Data availability

To ensure availability to the whole project team, the data files were downloaded from the kaggle website stated in the project brief and stored on the project storage location at Amazon S3:

- 2020\_spotify.csv (252.31 KB)
- brazil\_data.csv (1.29 MB)
- world\_data.csv (1.28 MB)

# Process workflow / pipeline

The project team agreed on the following workflow to conduct the required tasks.

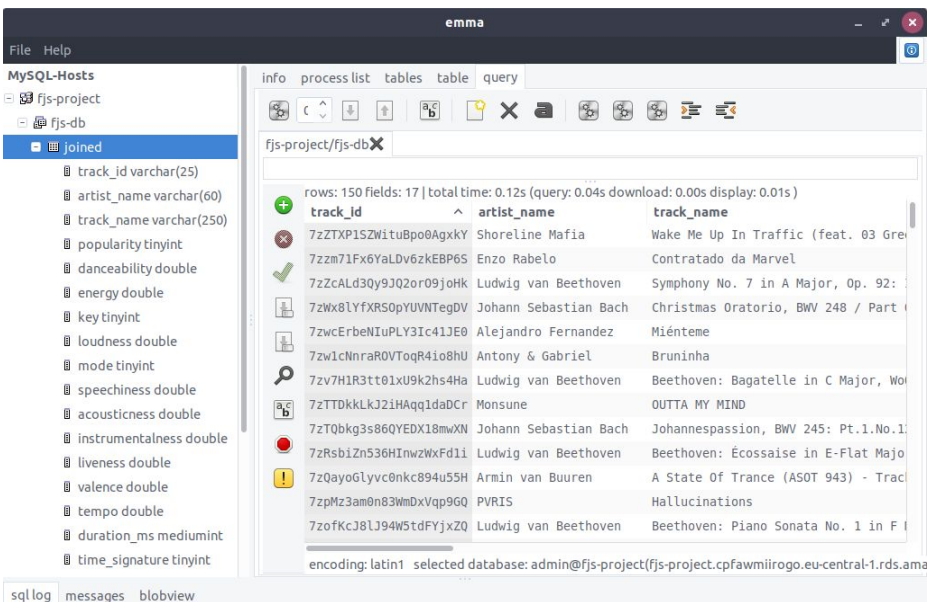
## Data preprocessing

- We received music title data in 3 CSV-file datasets:  
'world': 9,320 rows, 'brazil': 9,239 rows, '2020': 1,742 rows.
- These CSV-files were stored on our S3-bucket and retrieved from there by our Sagemaker notebooks.
- After concatenating all rows (= 20,301 rows), we removed exact (-4481 = 15,820 rows) and high-similarity duplicates (-194 = 15,626 rows).

```
# Now some more complicated duplicates - same artist, track_name and duration  
subset = 'artist_name track_name duration_ms'.split()  
df_joined[df_joined.duplicated(subset=subset, keep=False)].sort_values(subset)
```

	artist_name	track_name	track_id	popularity	danceability	energy	key	loudness	mode	speechiness	acousticness	instrumentalness
4268	As I Lay Dying	Blinded	5xa5C8TmCuF1q8cBRutiMY	62	0.339	0.9720	8	-5.260	1	0.1850	0.000006	0.
8774	As I Lay Dying	Blinded	2HdjEa5BP2VACT1velDTik	56	0.332	0.9720	8	-5.258	1	0.1980	0.000006	0.
3911	As I Lay Dying	My Own Grave	6bDoeNLCA32SYhTzk5w5y	61	0.475	0.9940	5	-4.978	0	0.2500	0.000113	0.
6751	As I Lay Dying	My Own Grave	0CcQWuAEJC93K8cBMbAigl	57	0.477	0.9940	5	-4.953	0	0.2410	0.000115	0.
3110	Bastille	Admit Defeat	1OHAfGgB1Jatvto7HvUN4L	56	0.653	0.6410	5	-7.007	1	0.0475	0.155000	0.
1952	Bastille	Admit Defeat	4qcnL8k4i8Ynw7kAPgVoD6	54	0.651	0.6400	5	-7.019	1	0.0460	0.149000	0.
6764	Cigarettes After Sex	Cry	0r0zdaQ9S3fouDnvJ25pwl	63	0.408	0.3970	7	-10.452	1	0.0279	0.756000	0.
8581	Cigarettes After Sex	Cry	0Qr61NXIlyAeQaADO5xn3rl	53	0.409	0.3990	7	-10.456	1	0.0275	0.763000	0.

- The joined dataset was written to a CSV-file on our S3 bucket and to a table on an RDS database instance.



## Exploratory Data Analysis (EDA)

- Analysis of the data completeness showed that in all 14 feature columns, there were no null values.



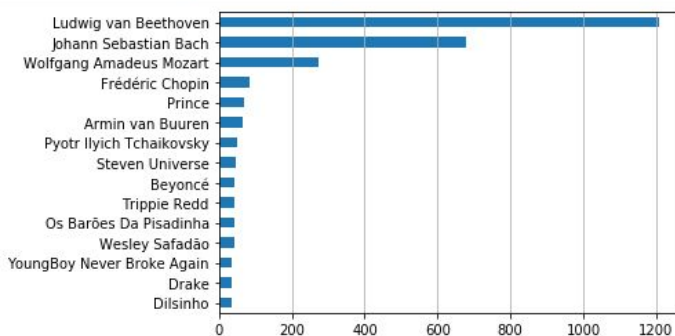
```
df_joined.describe().rename_axis('column').T.astype({'count':int})
```

	column	count	mean	std	min	25%	50%	75%	max
	popularity	15626	49.098234	24.273116	0.00000	43.00000	58.000000	65.000000	100.000
	danceability	15626	0.594562	0.195956	0.00000	0.45800	0.622000	0.745000	0.983
	energy	15626	0.534537	0.271949	0.00002	0.33700	0.588000	0.750000	1.000
	key	15626	5.183348	3.594108	0.00000	2.00000	5.000000	8.000000	11.000
	loudness	15626	-10.063025	7.372410	-45.13600	-11.83775	-7.164500	-5.221000	2.036
	mode	15626	0.633239	0.481936	0.00000	0.00000	1.000000	1.000000	1.000
	speechiness	15626	0.109787	0.118288	0.00000	0.03990	0.056200	0.125000	0.951
	acousticness	15626	0.417139	0.362938	0.00000	0.07410	0.309000	0.783000	0.996
	instrumentalness	15626	0.167178	0.333673	0.00000	0.00000	0.000003	0.014075	0.998
	liveness	15626	0.200374	0.184837	0.00000	0.09760	0.123000	0.227000	1.000
	valence	15626	0.475397	0.250612	0.00000	0.27400	0.470500	0.671000	0.989
	tempo	15626	119.823206	30.937673	0.00000	95.01500	120.143500	140.030750	235.998
	duration_ms	15626	204291.707283	96346.332970	12564.00000	160112.00000	191793.000000	226003.000000	1493227.000
	time_signature	15626	3.884423	0.502236	0.00000	4.00000	4.000000	4.000000	5.000

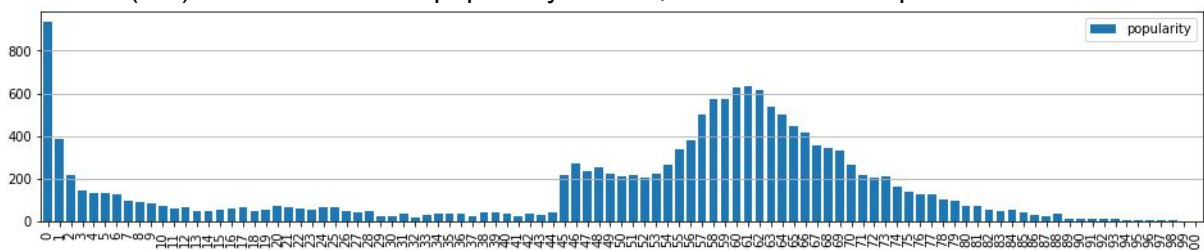
- Analysis of the “artist” and “popularity” features showed a very strong bias to classical music: of 15,626 rows, 1206 alone belong to Ludwig van Beethoven, a further 679 to Johann Sebastian Bach and another 275 to Wolfgang Amadeus Mozart.

These three composers already cover roughly 14% of all titles!

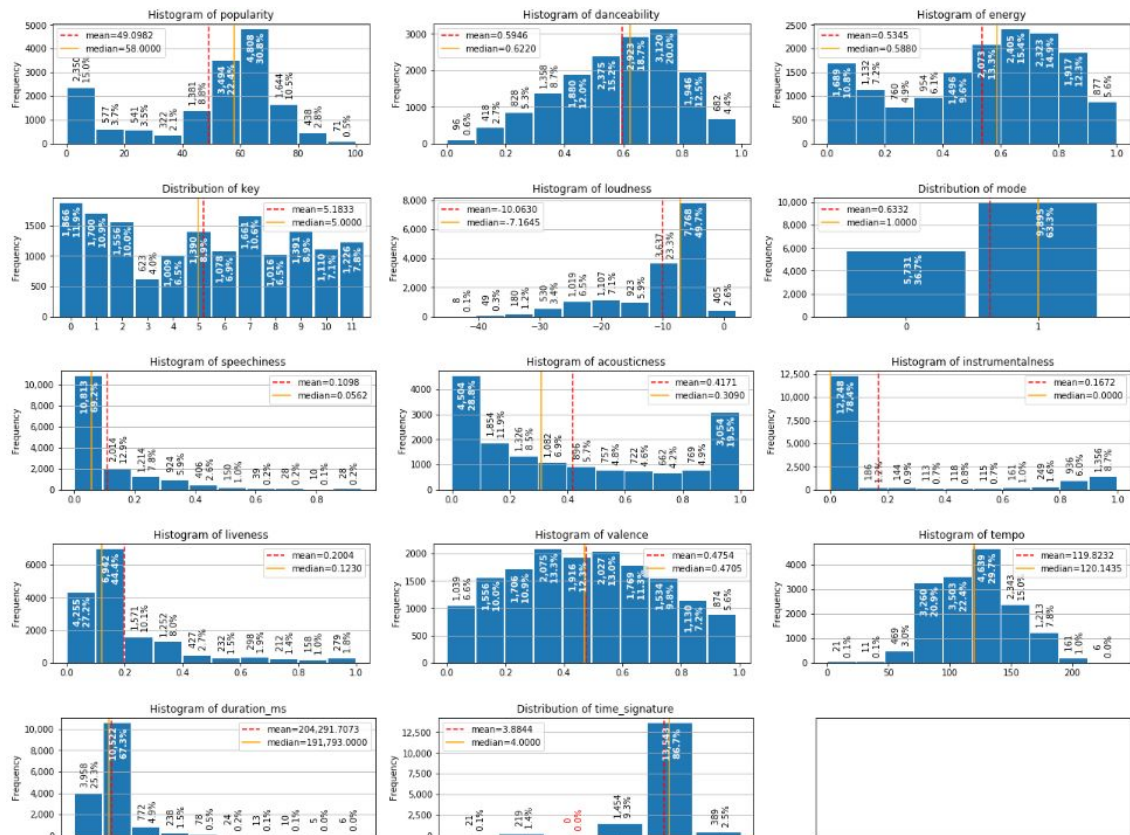
```
df_joined['artist_name'].value_counts()[15::-1].plot.barh().grid(axis='x')
```



- 938 rows (6%) of all titles have a popularity of zero, which seems suspect.



- Analysis of value ranges showed which features need scaling for use in clustering algorithms.



- Also, we looked for outliers per feature and quantified their volume in the data.

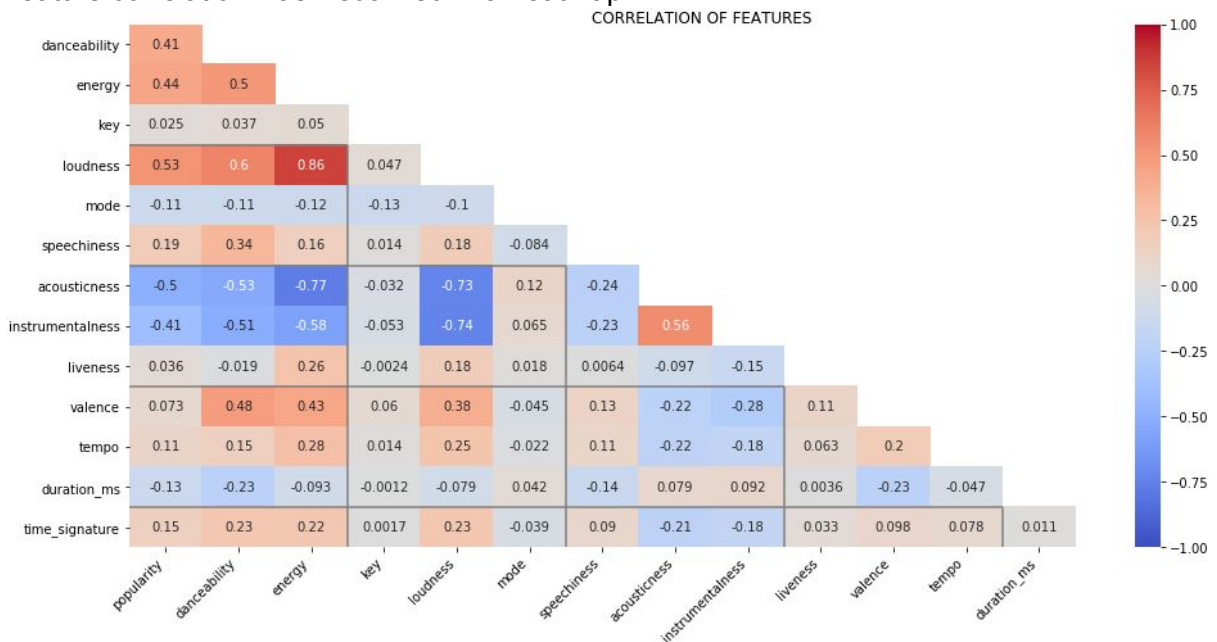
```
print(f'Total rows of data in df_joined dataframe: {len(df_joined):,}')

df_trimmed = df_joined.query('!(loudness < -35) or (loudness > 0)
or (speechiness > 0.65) or (tempo < 1) or (duration_ms > 430000)
or (time_signature < 1)')

draw_hist_plots('df_trimmed')
```

Total rows of data in df\_joined dataframe: 15,626  
Total rows of data in df\_trimmed dataframe: 609

- Some playful comparisons of titles from different artists were performed.
- Feature correlation was visualized in a heatmap:



## Clustering Algorithms

The project presents an Unsupervised Clustering problem, for this the following algorithms are recommended:

- K-means
- Hierarchical
- DBScan

The DBScan algorithm is recommended, when the K-means and hierarchical reach their limits due to having to handle noise and generating bad results with finding clusters of nonspherical shape. (see Amit Shreiber, A Practical Guide to DBSCAN Method)

## Distributed Cluster-Computing Frameworks

### AWS Cloud

#### AWS Sagemaker Standalone Notebook

##### **Initial Validation using Kmeans model from sklearn**

To enable a direct comparison between the distributed and standalone computing frameworks the team chose to implement the primary clustering model also on an AWS Sagemaker Standalone Notebook. The implementation (K-Means-SKLearn\_jgo.ipynb) required a standard Jupyter notebook and utilized a sklearn kmeans model. The algorithm was set with a parameter of n\_clusters=24 (the same setting as on the distributed-computing experiment).

The notebook completed the computation with the following processing times captured:

user 2.65 s, sys: 228 ms, total: 2.88 s Wall time: 2.73

This shows that the project dataset does not require large computational resources.

Therefore it was decided to implement other models on the standalone notebook setup.

##### **DBScan model from sklearn**

DBScan was selected as the second clustering model and for this assessment the sklearn model was implemented. The initial implementation (DBScan\_jgo.ipynb) utilized the full dataset (joined.csv) and included basic EDA (heatmap of feature correlation) and data preprocessing (dropping of text based columns and scaling of numeric values) steps for transparency (the detailed process was being developed at this time).

Before the data could be clustered it was necessary to assess which hyperparameter set (eps and min\_samples) would best suit the task.

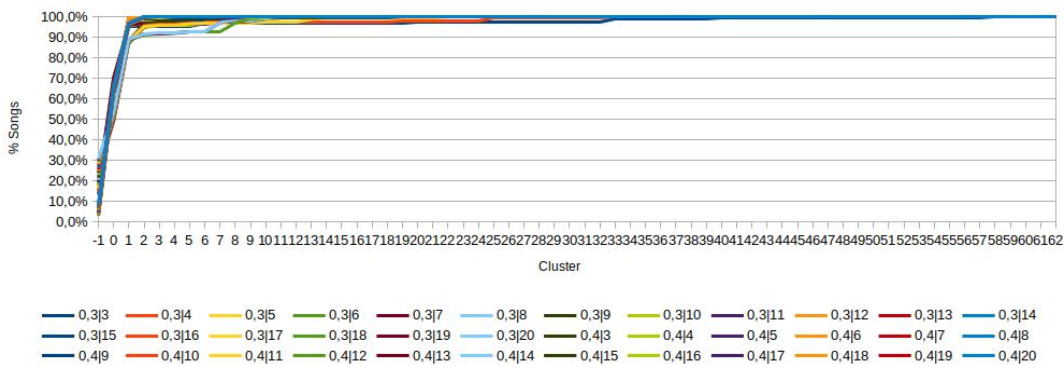
This was conducted by using a nested for-loop with ranges of eps and min\_samples values.

Unfortunately the activity led to frequent outages due to Gateway timeout and dead kernel errors between the Jupyter Notebook and the AWS Sagemaker server backend.

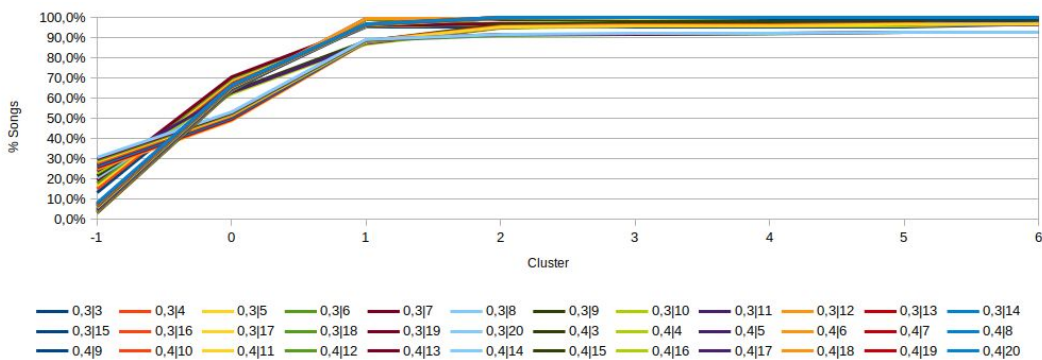
Reassessment of the source data during this phase gave rise to the insight that the zero popularity values were an anomaly and should be removed to provide a “train” dataset and a potential “test” dataset for future prediction efforts should they become relevant.

In a second iteration of the notebook the hyperparameter assessment was split into multiple cells to reduce individual processing time of individual cells and to avoid a timeout error that occurs after a cell “walltime” of approximately 10 minutes.

This modification resulted in completed hyperparameter assessment and the resultant cluster dynamics were then transferred to a diagram showing the cumulative percentage of total songs for the relevant clusters by hyperparameter set (shown in the legend as eps|min\_samples) for evaluation.



Here the same diagram but zoomed into the left hand section to show the relevant section.



Following guidance on selecting parameters for DBScan:

"... According to a research made in 2017 by Schubert, Sander, et al, the desirable amount of noise will usually be between 1% and 30%.

Another insight from that research is that if one of the clusters contains many (20%-50%) points of the dataset, it indicates that you should choose a smaller value for  $\epsilon$  or to try another clustering method. ..." (Amit Shreiber, A Practical Guide to DBSCAN Method)

Two sets were selected: eps=0.3, min\_samples=5 and eps=0.3, min\_samples=18

The first was deployed in a separate notebook (DBScan\_deployment\_jgo.ipynb) which generated a first clustering.

The notebook was modified to take the insight on excluding zero "popularity" values. This was performed in the data preprocessing section. This notebook is named DBScan\_no\_pop0\_jgo.ipynb.

The results were more meaningful and as a final modification toward a faster and less error prone version, the individual pipeline process modules were split into the following separate notebooks:

Data Pre-Processing (Data\_Wrangle\_no\_pop0\_jgo.ipynb)

Steps include:

- loading source data from RDS (15626 rows × 17 columns)
- Split into:
  - df\_train: (14688, 17)
  - df\_test : (938, 17)



Encoding of “train” data (Encoding\_Train\_dataset\_no\_pop0\_jgo.ipynb)

Steps include:

- dropping non numeric columns ('artist\_name' and 'track\_name')
- scaling of numeric columns

Deployment of selected hyperparameter

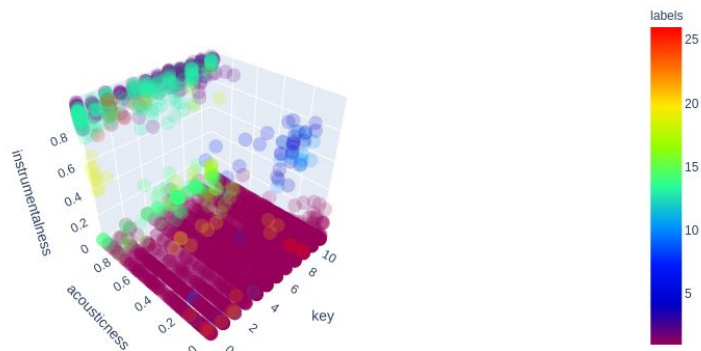
sets 0.3 | 5eps=0.3, min\_samples=5 and 0.3 | 18 eps=0.3, min\_samples=18

(DBScan\_deployment\_no\_pop0\_jgo.ipynb and DBScan\_deployment\_no\_pop0\_jgo-set2.ipynb)

Steps include:

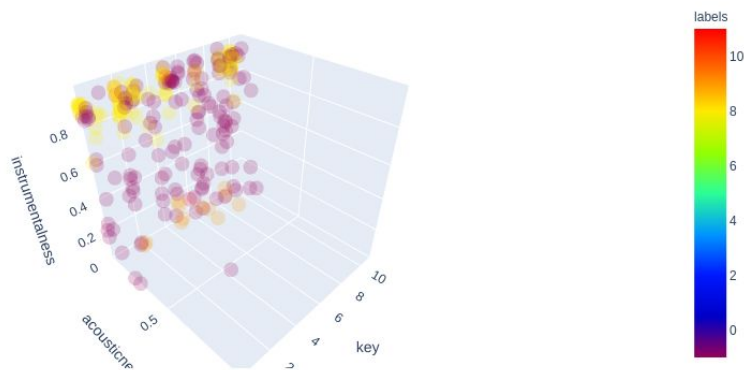
- Clustering to resulting clusters
  - Visualization of Clusters and individual artists using 3d scatter plots in Plotly
- e.g. Clusters  $\geq 0$

Cluster diagram for all Clusters  $\geq 1$



e.g. WAM Tracks

Cluster diagram for artist\_name Wolfgang Amadeus Mozart



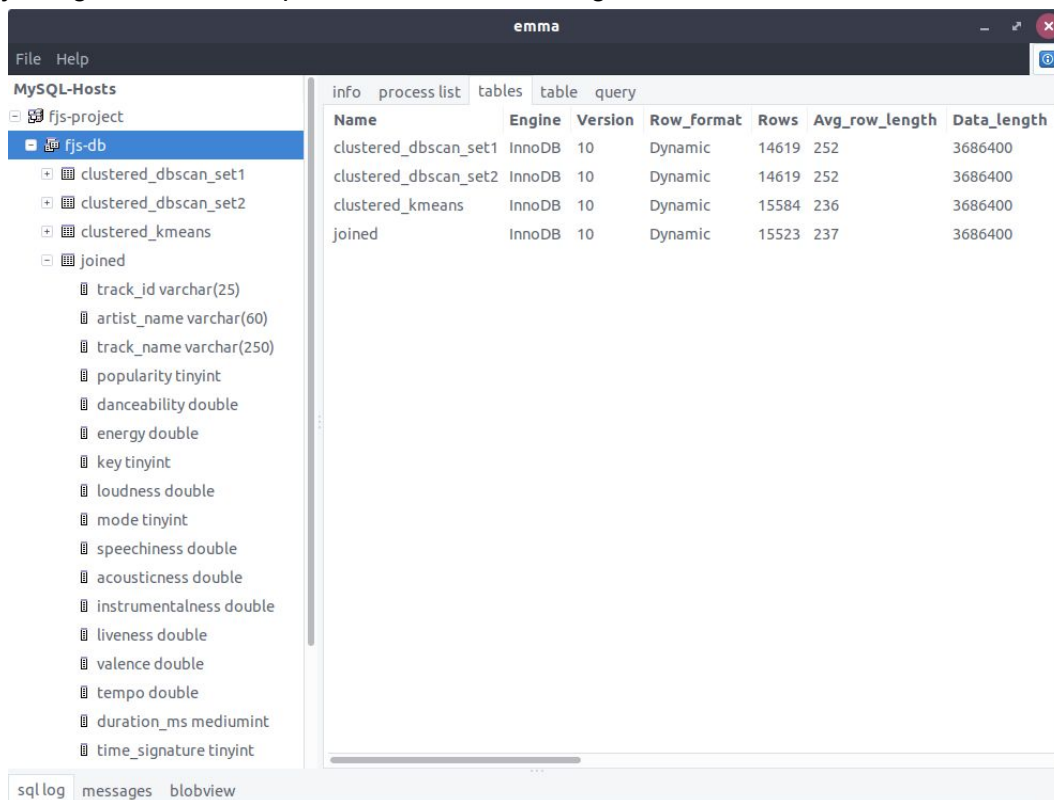
- Storing resultant DataFrame to RDS

AWS Sagemaker Cluster

## Apache Spark on Databricks

- Using the builtin PySpark ml-module, simple kMeans clustering was performed.  
The results did not differ significantly from the finding using kMeans in Sagemaker.

- The clustering assignment was written back to the RDS instance as per our project requirement, joining the same dumps from the DBSCAN algorithms:



Name	Engine	Version	Row_format	Rows	Avg_row_length	Data_length
clustered_dbscan_set1	InnoDB	10	Dynamic	14619	252	3686400
clustered_dbscan_set2	InnoDB	10	Dynamic	14619	252	3686400
clustered_kmeans	InnoDB	10	Dynamic	15584	236	3686400
joined	InnoDB	10	Dynamic	15523	237	3686400

## Learnings / Insights

### Clustering and Recommendation

- From the project description, we understood the primary goal to be clustering of the titles, which we performed successfully (albeit the available features produced sometimes “weird” cluster companions).
- As a secondary goal, we understood a kind of recommendation based on the clustering should be evaluated. This target can fairly squarely be regarded as failed, because the “weird” cluster composition mentioned above makes good recommendations extremely unlikely.  
An example: we found no less than 36 titles belong to one recording of the opera “Leonore”. Due to large variations in feature assignment (in an opera, you have arias, a-cappellas, choruses, interludes, you have many different tempi and pieces in various keys, happy and sad moods...), these 36 titles were clustered by kMeans into no less than 23 different clusters!  
Obviously, a “recommendation” for even the next piece within the opera based on this cluster assignment failed miserably. Instead, probable outliers in the clustering joined pieces from Ludwig van Beethoven and Kanye West in the same cluster!

### Performance

Identical kMeans-tasks were run in three different environments:

- AWS Sagemaker Cluster using an Sagemaker-internal kMeans method
- Standalone execution within an AWS Sagemaker notebook using sklearn module
- Within a DataBricks workspace

The comparison of processing time was revealing:

```
2020-12-03 21:11:14 Uploading - Uploading generated training model
2020-12-03 21:11:14 Completed - Training job completed
Training seconds: 46
Billable seconds: 46
CPU times: user 779 ms, sys: 34.6 ms, total: 813 ms
Wall time: 3min 12s
```

Sagemaker Cluster-computing:

sklearn Standalone notebook: CPU times: user 2.65 s, sys: 228 ms, total: 2.88 s Wall time: 2.73 s

For the small data volume we worked on in this project, standalone sklearn was fastest by far.

The overhead of starting up the Sagemaker cluster, or the distributed Hadoop clustering in Databricks was much too large for the execution time to matter at all.

Learning: Use distributed cluster-computing only if the data requires it.

The time overhead is considerable.

## Potential next steps

- Remembering strong bias to classical composers in the source data, a more balanced data source should be obtained.
- Within the existing data, a noticeable outlier-range of popularity close to or equal to zero was observed. Possibly, based on the other features, a prediction of popularity could improve the distribution of that feature and thus also improve clustering results.
- Use the text columns (artist, track\_name) as features for the clustering, for example using word encoding

## Distribution of tasks

The individuals project members executed the following project tasks:

- Falk Lutz
  - kMeans clustering in SageMaker and DataBricks
- Julian Godley
  - Source data and infrastructure assurance
  - Clustering methods on Standalone Notebooks: DBScan, kMeans
- Simon Harston
  - CSV-file loading, concatenation, data preprocessing, exploratory data analysis
  - storage of source data and clustered data to S3 bucket and RDS database
  - setup of DataBricks environment and loading of data from RDS instance
  - debugging of kMeans clustering on DataBricks including
    - heatmap-style display of feature-to-cluster weights
    - thoughts on “recommendations” based on cluster assignment

## List of links

### Kaggle Challenge and other websites dealing with the idea

- Challenge site  
<https://www.kaggle.com/rafaelnduarte/spotify-data-with-audio-features>
- “Clustering the Most Listened to Songs of the 2010s Using Spotify Data.”  
<https://medium.com/analytics-vidhya/clustering-most-listened-songs-of-the-2010s-using-spotify-data-8e25e8b082ce>

- Discovering similarities across my Spotify music using data, clustering and visualization  
<https://towardsdatascience.com/discovering-similarities-across-my-spotify-music-using-data-clustering-and-visualization-52b58e6f547b>

## Spotify

- Glossary  
<https://artists.spotify.com/blog/glossary-of-music-terms-actual-music-terms>
- Documentation on audio analysis  
<https://developer.spotify.com/documentation/web-api/reference/tracks/get-audio-analysis/>
- Insight on Valence  
<https://community.spotify.com/t5/Content-Questions/Valence-as-a-measure-of-happiness/td-p/4385221>
- Insight on Mode & key  
<https://developer.spotify.com/documentation/web-api/reference/tracks/get-several-audio-features/>

## Clustering Algorithms

### Kmeans

...

### DBScan

Amit Shreiber, A Practical Guide to DBSCAN Method,

<https://towardsdatascience.com/a-practical-guide-to-dbscan-method-d4ec5ab2bc99>