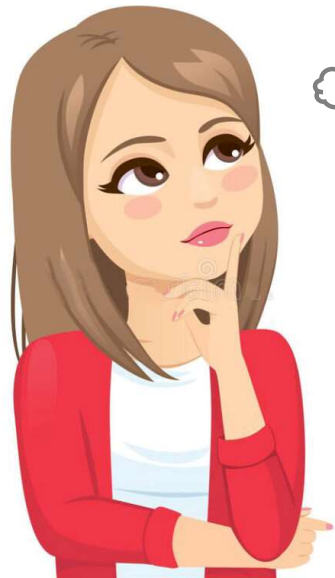




# Azure Rhythm

Cloud Miners  
SMU - 2023

The Idea...



# Presentation Journey

1. R Shiny Application - Haitie
2. ETL - Mai
3. Azure Database – Lijo
4. Azure Machine Learning – Todd

Guided by Lani



# Azure Rhythm

## R Shiny Application

Haitie Liu





# Azure Rhythm

R Shiny Application

<https://azrhythmwebapp.azurewebsites.net/>



# Extract, Transfer, Load (ETL)

Python

Mai Dang



# ETL



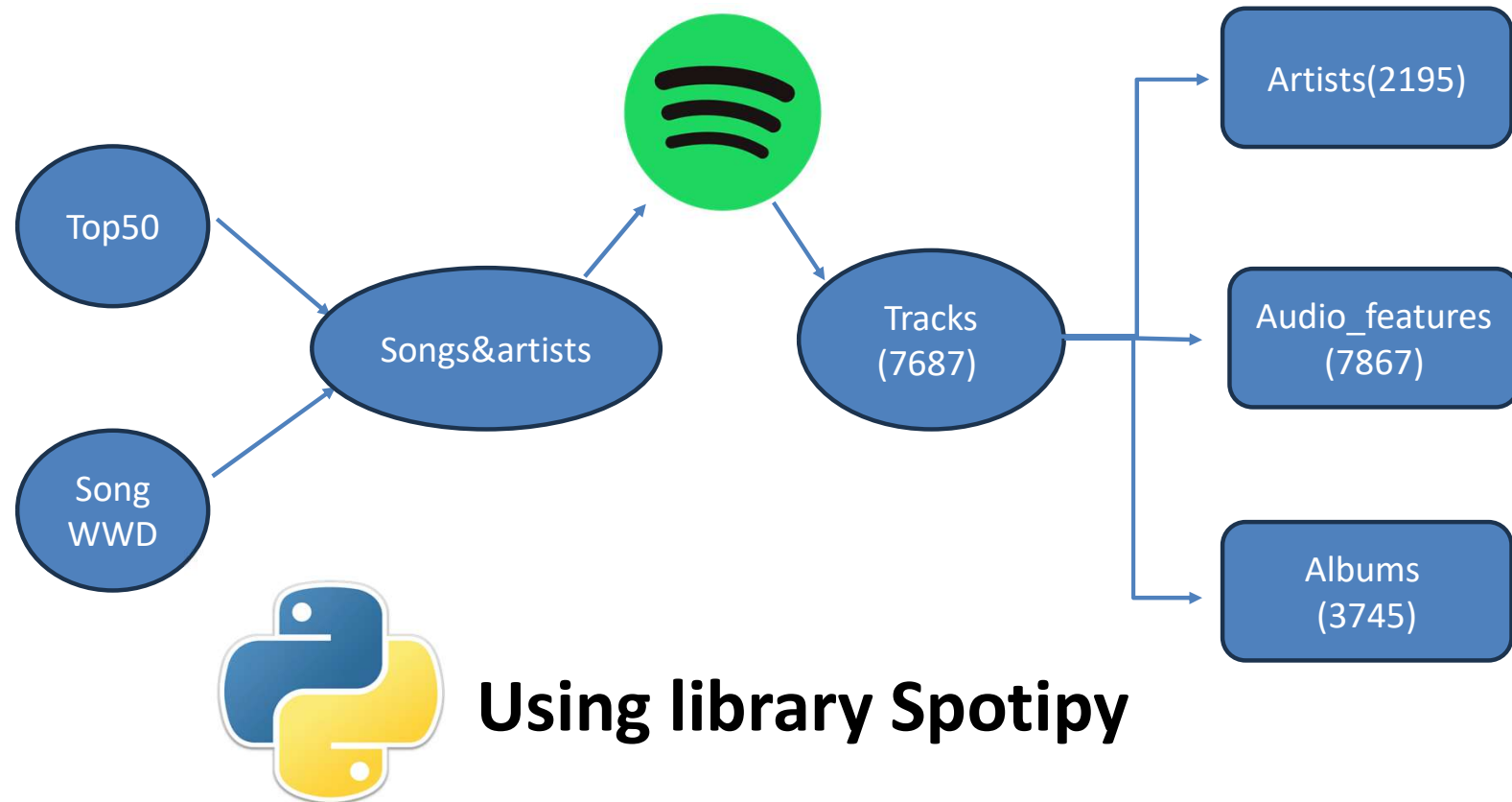
Extract

The diagram illustrates the ETL (Extract, Transform, Load) process. It features three rounded rectangular boxes arranged horizontally, each containing a step of the process. The 'Extract' box is orange, while 'Transform' and 'Load' are blue. These boxes are positioned over a large, light-blue arrow pointing to the right. The background of the slide is split: the top left has orange and yellow curved shapes, and the top right shows a blue-tinted image of electronic equipment.

Transform

Load

# Extract





# Extract

## Roadblocks

- Spotify API limits: rate limit, token limit, etc.
- Duplicated IDs
- Inconsistent artist names in 2 datasets

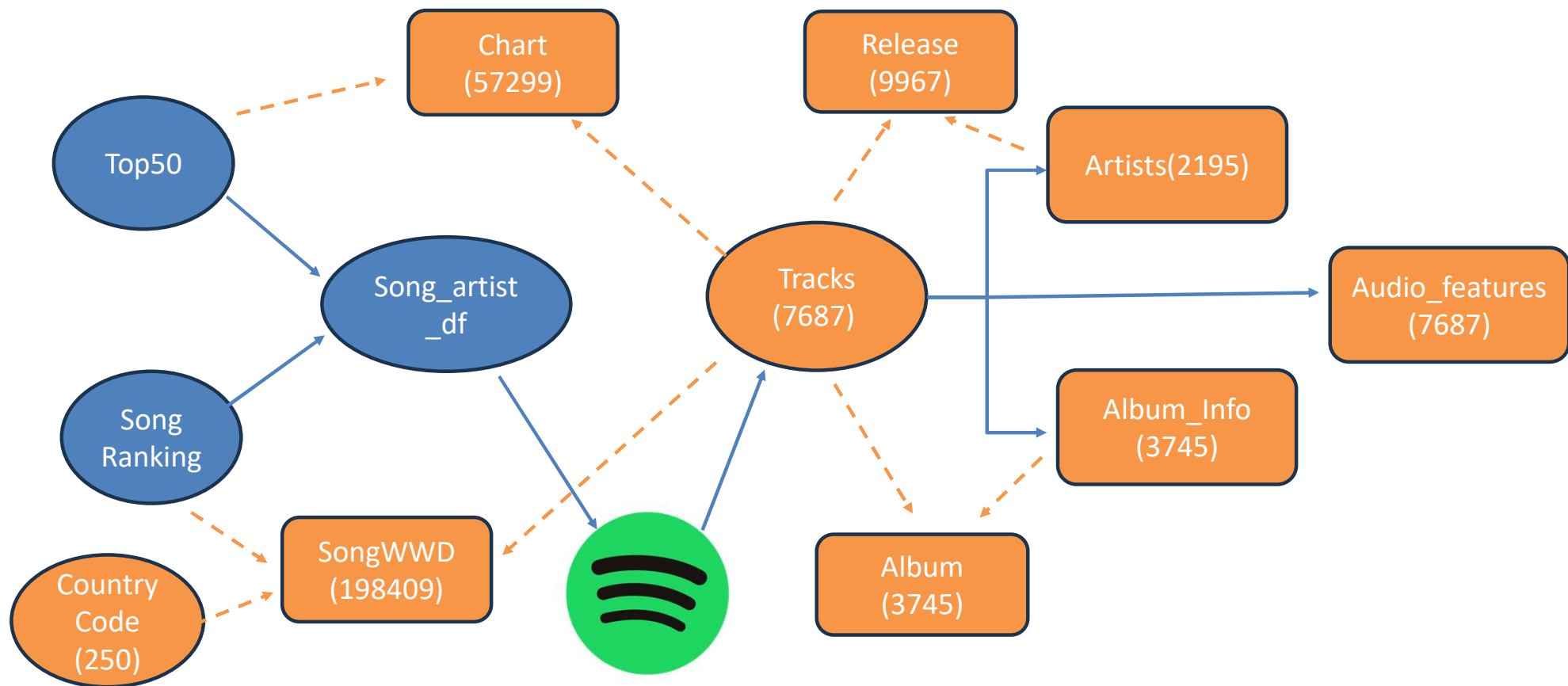
## Solutions

- Client Credential authentication  
Break between loops  
Response Counter.
- Check and remove duplicates.
- Remove symbols and characters before normalizing the artist names.

# ETL



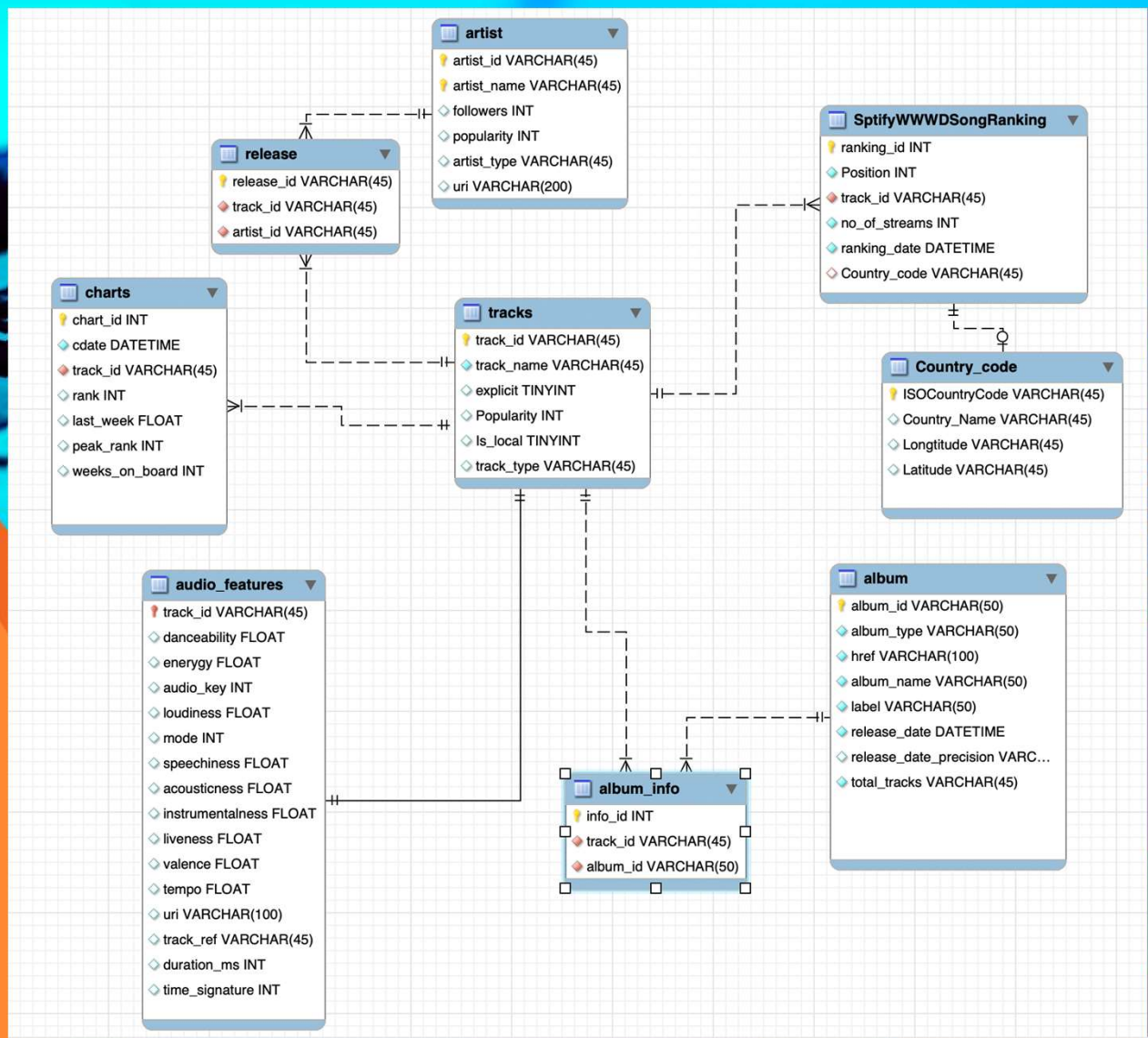
# Transform



# ETL



# ERD

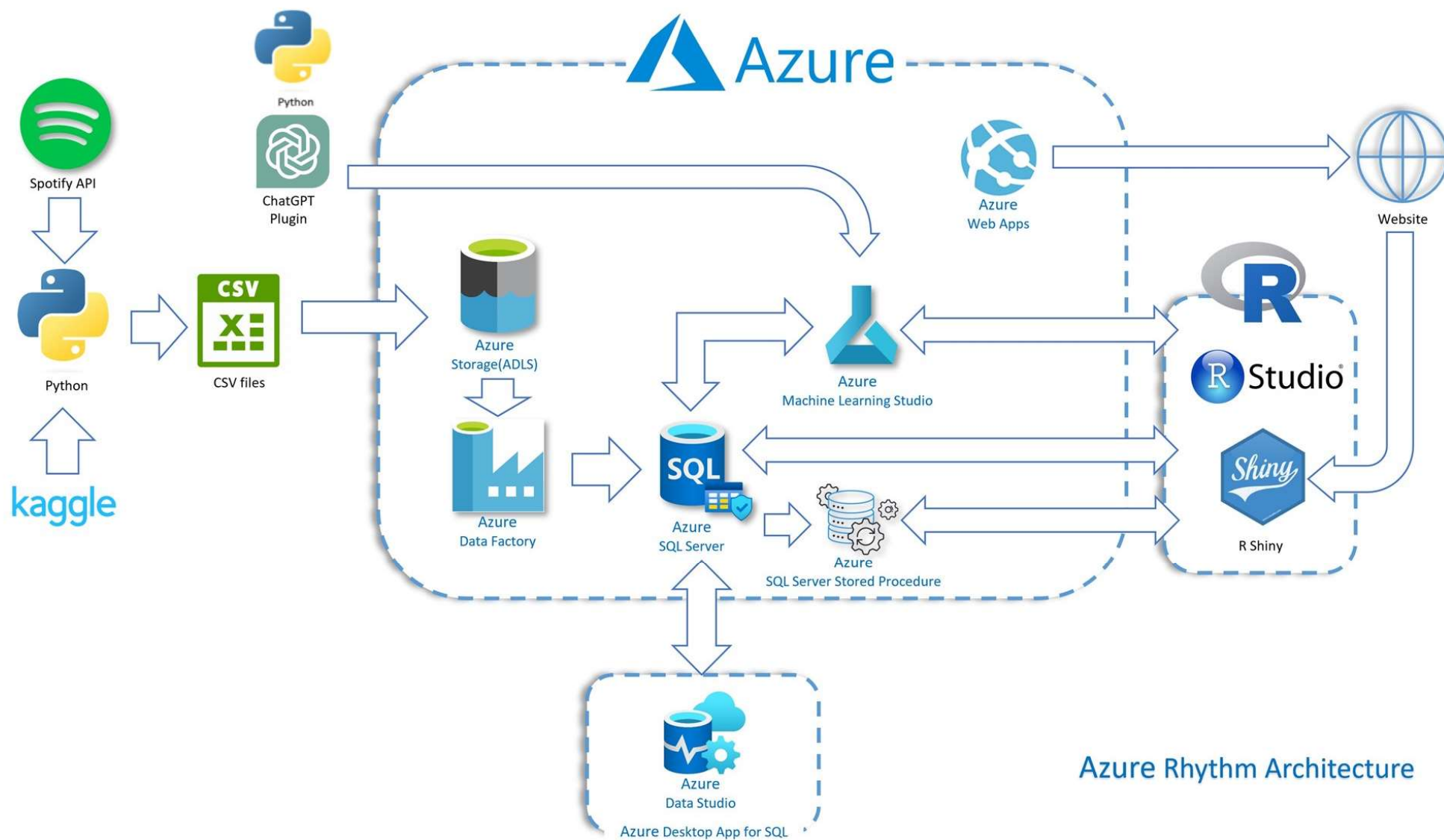




# Azure Database

Lijo Jacob





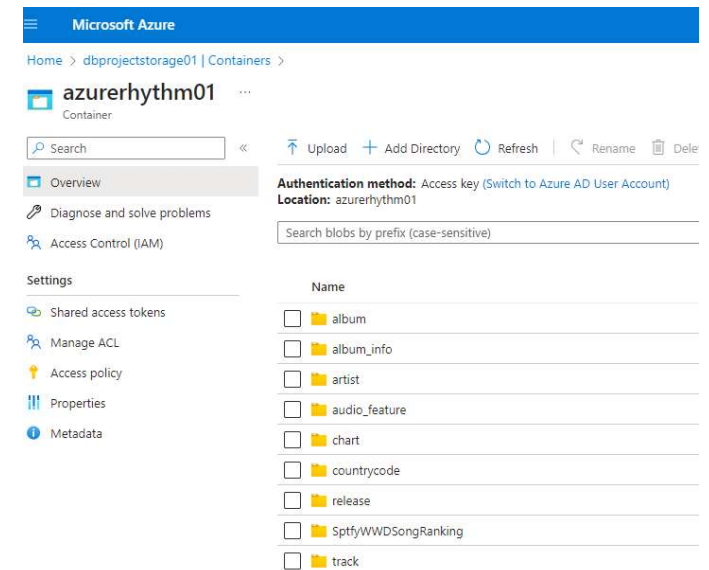
# Azure Setup & Data Migration to Azure

- Setup Azure student account for each team member.
- Created resources under one of our subscriptions and assigned “Contributor” role to each team member.
- Setup a storage account and created a container
- Data stored in csv files were migrated to Azure data lake storage.



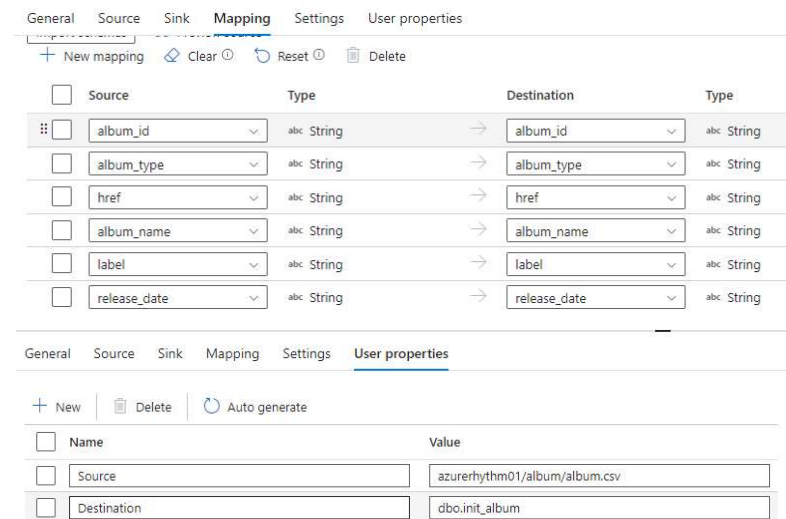
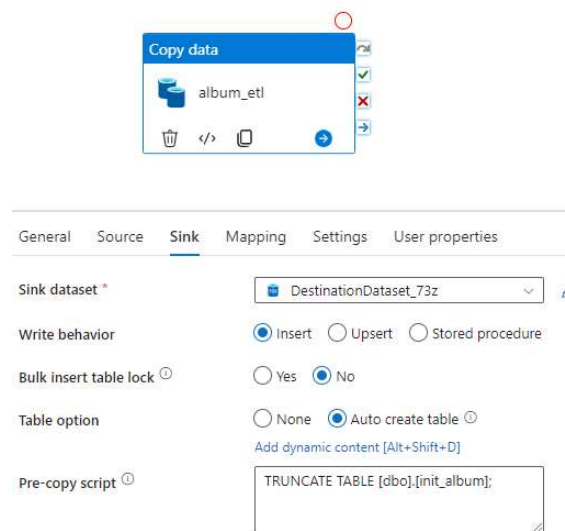
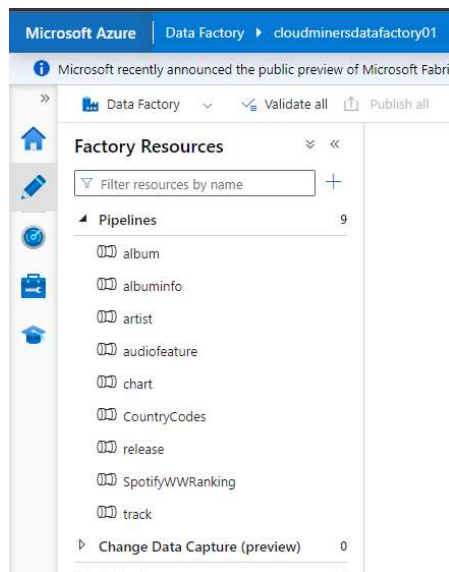
Microsoft Azure  
Data Lake

\*



# Data Ingestion to Azure SQL

- Data stored Azure data lake storage is ingested into Azure SQL using Azure Data Factory



# Azure SQL Setup



- Created a sql server and sql database
- Setup firewall rules to allow our machines to access the sql server using Azure data studio and R.

Microsoft Azure

Home > Cloud Miners | Resource groups > AzureRhythm > cloudminorsqlsverdb >

dbprojectsqlsvr (cloudminorsqlsverdb/dbprojectsqlsvr) SQL database

Search resources, services, and docs (G+/I)

Copy Restore Export Set server firewall Delete Connect with... Feedback

Overview

Activity log

Tags

Diagnose and solve problems

Query editor (preview)

Settings

Compute + storage

Connection strings

Properties

Locks

Data management

Replicas

Sync to other databases

Integrations

Azure Synapse Link

Stream analytics (preview)

Add Azure Search

Essentials

Resource group (move): [AzureRhythm](#)

Status: Online

Location: East US

Subscription (move): [Cloud Miners](#)

Subscription ID: 9dfed841-3609-4e5d-ab84-15399bcc55aa

Tags (edit): [Add tags](#)

Getting started **Monitoring** Properties Features Notifications (0) Integrations Tutorials

Database data storage

Review the below metrics and monitor your applications and infrastructure.

0.04% Used

Used space: 104.38 MB

Remaining space: 249.9 GB

Allocated space: 240 MB

Max storage: 250 GB

Server name: [cloudminorsqlsverdb.database.windows.net](#)

Elastic pool: [No elastic pool](#)

Connection strings: [Show database connection strings](#)

Pricing tier: Standard S0: 10 DTUs

Earliest restore point: 2023-07-23 03:29 UTC

Firewall rules

Allow certain public internet IP addresses to access your resource. [Learn more](#)

+ Add your client IPv4 address (██████████) + Add a firewall rule

Rule name	Start IPv4 address	End IPv4 address	
ClientIPAddress_2023-6-4_12-12-55	██████████	██████████	🗑️
ClientIPAddress_2023-6-4_18-48-35	██████████	██████████	🗑️
ClientIPAddress_2023-6-4_18-48-37	██████████	██████████	🗑️
ClientIPAddress_2023-6-4_18-48-38	██████████	██████████	🗑️
ClientIPAddress_2023-6-4_18-49-57	██████████	██████████	🗑️



# Relational Database Tables using Azure SQL

- For data ingested using Azure Data Factory, the landing tables are named as “init\_<file\_name>”.
- Relational tables are created with key constraints and appropriate data types as per the ER Diagram.



```
> dbo.init_album
> dbo.init_album_info
> dbo.init_artist
> dbo.init_audio_feature
> dbo.init_chart
> dbo.init_country_codes_Lat_long
> dbo.init_release
> dbo.init_SptfyWWDRanking
> dbo.init_track
```

```
> dbo.album
> dbo.album_info
> dbo.artist
> dbo.artists_v2
> dbo.audio_feature
> dbo.chart
> dbo.Country_code
> dbo.release
> dbo.SptfyWWDSongRanking
> dbo.track
```

# Relational Table Creation and Data Insertion

- Table are created with Primary Keys and Foreign keys using create table statements.
- Data is inserted into the table from corresponding landing table with appropriate foreign key table lookup.

```
CREATE TABLE [dbo].[audio_feature](  
    [af_track_id] [nvarchar](50) PRIMARY KEY NOT NULL,  
    [Danceability] [float] NOT NULL,  
    [Energy] [float] NOT NULL,  
    [Audio_key] [tinyint] NOT NULL,  
    [Loudness] [float] NOT NULL,  
    [Mode] [tinyint] NOT NULL,  
    [Speechiness] [float] NOT NULL,  
    [Acousticness] [float] NOT NULL,  
    [Instrumentalness] [float] NOT NULL,  
    [Liveness] [float] NOT NULL,  
    [Valence] [float] NOT NULL,  
    [Tempo] [float] NOT NULL,  
    [Uri] [nvarchar](150) NOT NULL,  
    [Track_herf] [nvarchar](150) NOT NULL,  
    [Duration_ms] [int] NOT NULL,  
    [Time_signature] [tinyint] NOT NULL,  
    FOREIGN KEY (af_track_id) REFERENCES [dbo].[track] (track_id)  
);
```

```
insert into dbo.audio_feature select * from dbo.init_audio_feature where af_track_id in (select track_id from dbo.track);
```



# MACHINE LEARNING THE FORMULA FOR SUCCESS!

Todd Garner





# Why Azure Machine Learning?

# Introduction and Problem Statement

- 11 independent variables from Spotify, one response variable (popularity).

popularity	followers	Danceability	Energy	Loudness	Mode	Speechiness	Acousticness	Instrumentalness	Liveness	Valence	Tempo
100	78485332	0.519	0.527	-7.673	1	0.0274	0.075	0	0.132	0.267	78.915
100	78485332	0.354	0.267	-13.69	1	0.0281	0.731	0.000402	0.0858	0.113	94.219
100	78485332	0.602	0.736	-5.778	1	0.0338	0.00196	4.57E-05	0.105	0.471	96.969
100	78485332	0.624	0.757	-2.94	1	0.0296	0.00265	1.87E-06	0.189	0.658	121.07
100	78485332	0.777	0.357	-6.942	1	0.0522	0.757	7.28E-06	0.108	0.172	139.883
100	78485332	0.472	0.701	-3.72	1	0.0279	0.091	0	0.23	0.304	147.854
100	78485332	0.392	0.574	-9.195	1	0.17	0.833	0.00179	0.145	0.529	81.112
100	78485332	0.636	0.402	-7.855	1	0.031	0.0494	0	0.107	0.208	125.952
100	78485332	0.58	0.491	-6.462	1	0.0251	0.575	0	0.121	0.425	76.009



# Methodology and Solution

## Derived Linear Regression Equation

**Popularity** =  $61.3331 + (3.93 \times 10^{-7}) \times \text{Followers} + (11.5315) \times \text{Danceability} - (0.0847) \times \text{Loudness} - (1.5251) \times \text{Mode} + (12.5477) \times \text{Speechiness} - (2.2467) \times \text{Acousticness} + (0.1258) \times \text{Liveness} - (6.2922) \times \text{Valence} + (0.0101) \times \text{Tempo}$

In this code snippet:

Regene

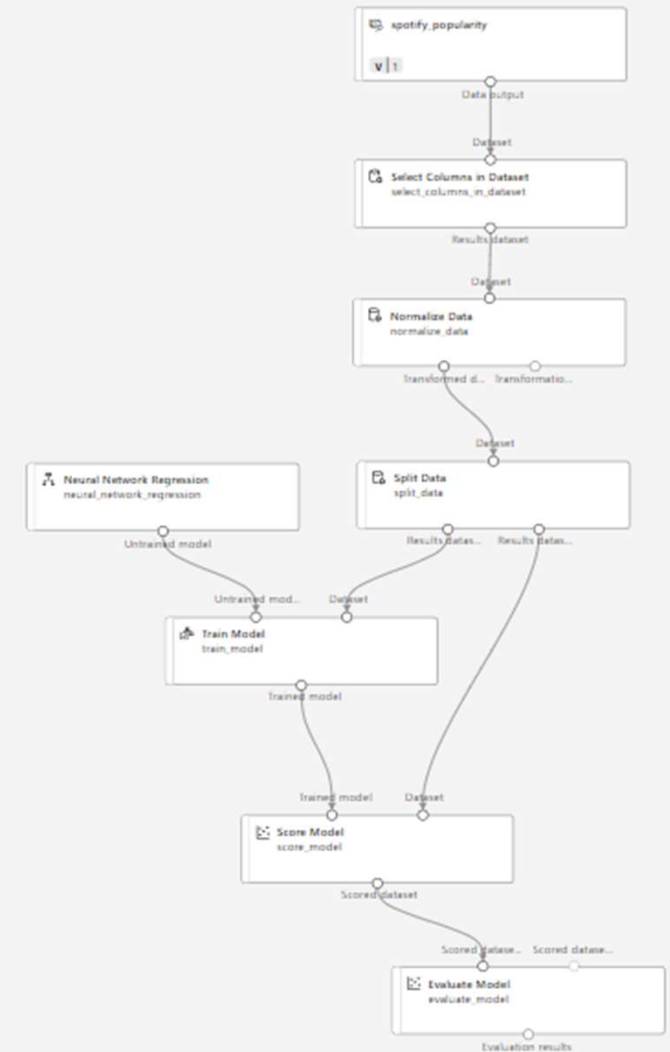
+ Send a message

ChatGPT may produce inaccurate information ab

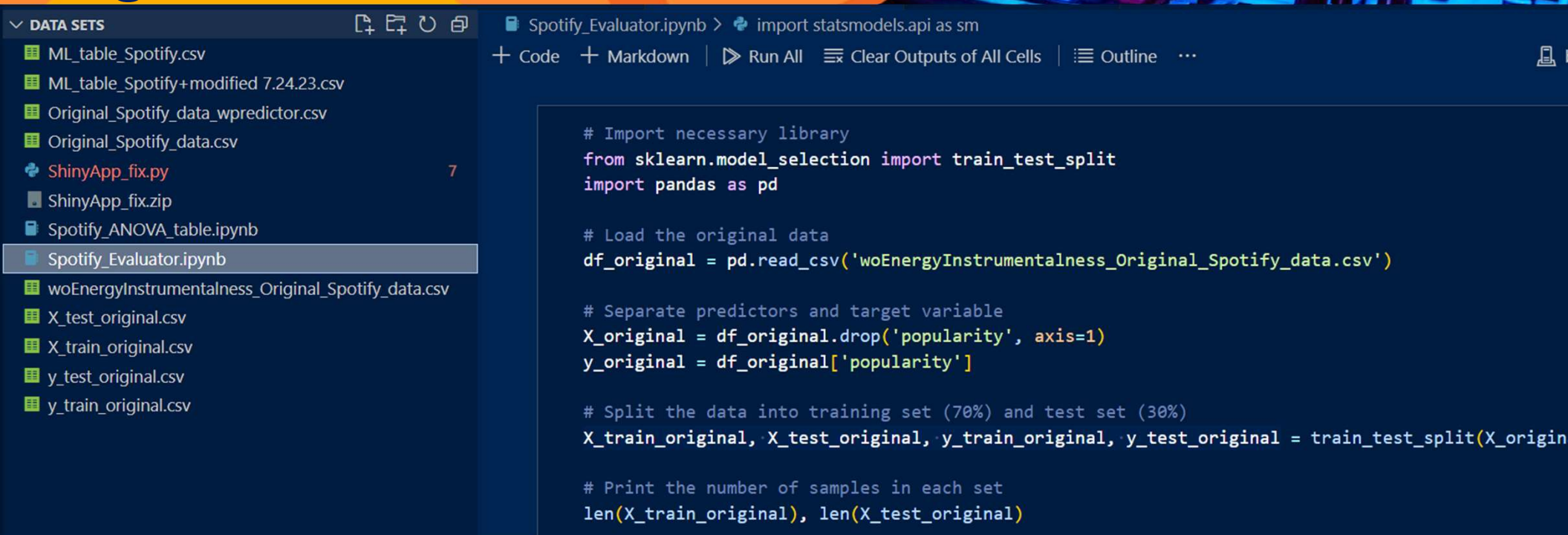
Data Component

95  

- ▶ Sample data (16)
- ▶ Data Transformation (19)
- ▶ Computer Vision (6)
- ▶ Model Scoring & Evaluation (6)
- ▶ Machine Learning Algorithms (19)
- ▶ Text Analytics (7)
- ▶ Python Language (2)
- ▶ Data Input and Output (3)
- ▶ Recommendation (5)
- ▶ R Language (1)
- ▶ Feature Selection (2)
- ▶ Anomaly Detection (2)
- ▶ Statistical Functions (1)
- ▶ Model Training (4)
- ▶ Web Service (2)



# VS Code – Python – Linear Regression



The screenshot displays the Visual Studio Code (VS Code) interface. On the left, the 'DATA SETS' file explorer shows a list of files: ML\_table\_Spotify.csv, ML\_table\_Spotify+modified 7.24.23.csv, Original\_Spotify\_data\_wpredictor.csv, Original\_Spotify\_data.csv, ShinyApp\_fix.py (with a '7' next to it), ShinyApp\_fix.zip, Spotify\_ANOVA\_table.ipynb, and Spotify\_Evaluator.ipynb (which is selected). The main editor area shows the code for Spotify\_Evaluator.ipynb. The code imports statsmodels.api as sm, then uses sklearn.model\_selection.train\_test\_split and pandas (pd) to load data from 'woEnergyInstrumentalness\_Original\_Spotify\_data.csv', drop the 'popularity' column to create X\_original, and split the data into training and testing sets (70% and 30% respectively). It also prints the number of samples in each set.

```
Spotify_Evaluator.ipynb > import statsmodels.api as sm
+ Code + Markdown | ▶ Run All | Clear Outputs of All Cells | Outline ...

# Import necessary library
from sklearn.model_selection import train_test_split
import pandas as pd

# Load the original data
df_original = pd.read_csv('woEnergyInstrumentalness_Original_Spotify_data.csv')

# Separate predictors and target variable
X_original = df_original.drop('popularity', axis=1)
y_original = df_original['popularity']

# Split the data into training set (70%) and test set (30%)
X_train_original, X_test_original, y_train_original, y_test_original = train_test_split(X_original, y_original, test_size=0.3, random_state=42)

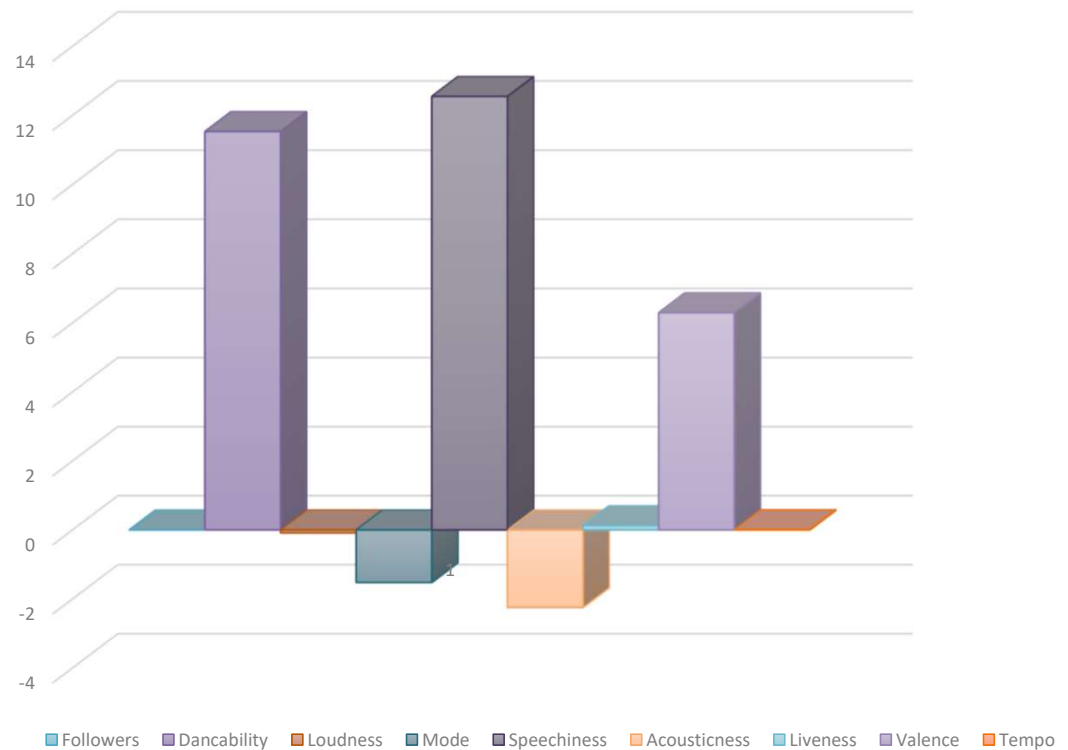
# Print the number of samples in each set
len(X_train_original), len(X_test_original)
```

Finally, I used VS Code and wrote a Python script and obtained a solution

# Results and Application

Popularity Equation parameter weights

Popularity	
3.93E-07	Followers
11.5315	Dancability
-0.0847	Loudness
-1.5251	Mode
12.5477	Speechiness
-2.2467	Acousticness
0.1258	Liveness
6.2922	Valence
0.0101	Tempo



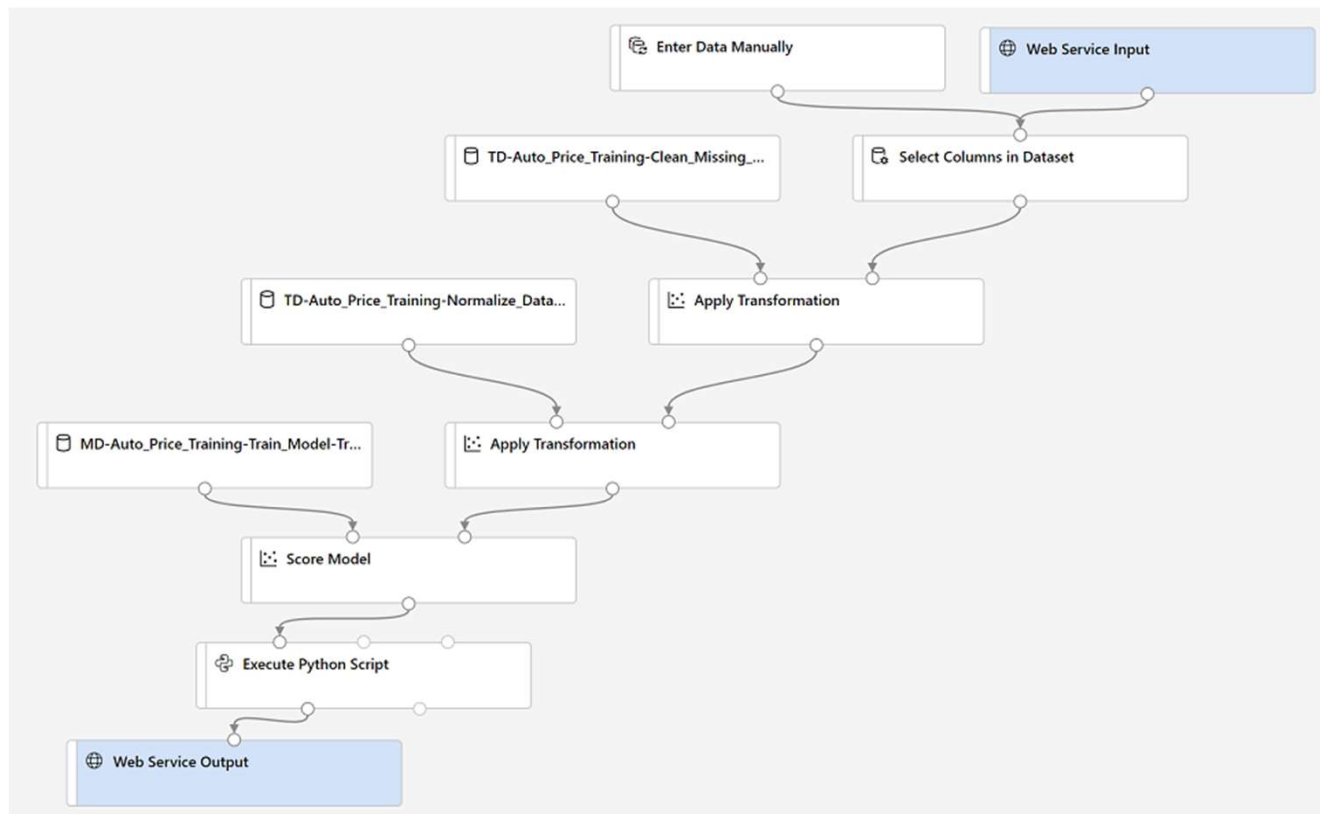
# Introduction and Problem Statement

- 9 independent variables from Spotify, one response variable (popularity), plus predicted popularity.

popularity	followers	Danceability	Loudness	Mode	Speechiness	Acousticness	Liveness	Valence	Tempo	popularity_predicted
100	78485332	0.519	-7.673	1	0.0274	0.075	0.132	0.267	78.915	96.14035357
100	78485332	0.354	-13.69	1	0.0281	0.731	0.0858	0.113	94.219	93.93886203
100	78485332	0.602	-5.778	1	0.0338	0.00196	0.105	0.471	96.969	95.61447613
100	78485332	0.624	-2.94	1	0.0296	0.00265	0.189	0.658	121.07	93.42389564
100	78485332	0.777	-6.942	1	0.0522	0.757	0.108	0.172	139.883	97.63948834
100	78485332	0.472	-3.72	1	0.0279	0.091	0.23	0.304	147.854	93.72877677
100	78485332	0.392	-9.195	1	0.17	0.833	0.145	0.529	81.112	93.50487737
100	78485332	0.636	-7.855	1	0.031	0.0494	0.107	0.208	125.952	97.27862442
100	78485332	0.58	-6.462	1	0.0251	0.575	0.121	0.425	76.009	94.87370596
100	78485332	0.316	-10.381	1	0.0488	0.878	0.0797	0.221	74.952	92.77302428
100	78485332	0.71	-6.965	1	0.0366	0.00164	0.0785	0.673	135.012	94.12101318



# Phase 2: Enhancements

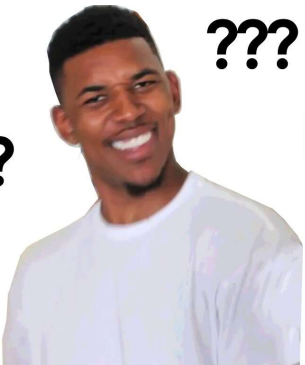




# Limitations

Is there a different way to load your data to the database?

Yes! We can use ODBC Driver 17 for Sql Server.





The End!

# Links

- [Azure Rhythm Application | Landing Page](#)
- [Rshiny App Break Down](#)
- [Coursera Machine Learning](#)



# Contacts

Haitie Liu - [haitiel@mail.smu.edu](mailto:haitiel@mail.smu.edu)

Mai Dang - [maid@mail.smu.edu](mailto:maid@mail.smu.edu)

Lijo Jacob - [lijoj@mail.smu.edu](mailto:lijoj@mail.smu.edu)

Todd Garner - [toddg@mail.smu.edu](mailto:toddg@mail.smu.edu)

Lani Lewis - [lanil@mail.smu.edu](mailto:lanil@mail.smu.edu)







Thank You!