# COMP4651 Final Project Report

# Music Recommendation with Cloud Integration

Fall 2024 | Group 5

KWONG, Hoi Wa
hwkwong@connect.ust.hk

SI TOU, Tin Nok
tnsitou@connect.ust.hk

TONG, Tsun Man
tmtong@connect.ust.hk

## Introduction

The purpose of our project is to leverage cloud computing to recommend the song to which a user might want to listen next, powered by machine learning (ML), based on various features related to the song, such as music genre, artists, and more. We have experimented with a multitude of methods, from model architecture to cloud computing deployment strategies, for optimal classification accuracy.

In this project, we make use of the KKBOX WSDM dataset [1], which features over 7 million entries of soundtrack playing history, together with attributes of the pieces, up till September 2017. Based on this dataset, we have engineered a PyTorch-based Feed-Forward Neural Network (FFNN) music genre classifier that consists of multiple fully connected (Dense) layers to identify correlations between the features of the song and its genre. We then deploy the classifier on an AWS EC2 instance for efficient data processing and model training. The code is available in our GitHub repository, where you can also find a demonstration video showcasing the project.

## Background

ML-based algorithms have demonstrated remarkable performance in music genre classification. Neural network architectures such as FFNN and Long Short-Term Memory (LSTM) are commonly implemented models. Frameworks like PyTorch provide the flexibility and scalability required to develop and train these models efficiently.

Amazon Web Services (AWS) Elastic Compute Cloud (EC2) is a prominent cloud computing service that provides resizable compute capacity, enabling developers to deploy, manage, and scale applications seamlessly. EC2 supports a wide range of instance types optimized for various workloads, including those requiring high computational power and storage capabilities essential for large-scale ML tasks.

## Related Work

Our work is inspired by relevant work in the open-source community, which aims at running neural networks on the cloud. One such example is Xing et al. [2] which extracts, from the mel-spectrogram of music tracks, the music latent factor vector using Weighted Matrix Factorization (WMF), which is subsequently passed to a Convolutional Neural Network (CNN) and LSTM model to distil latent spatial features of the source, combined with an embedding layer in parallel to predict users' preference for music based on the feature data from their playing history. While this solution is edge-cloud-based, which utilizes the resources of edge endpoints to operate the models, the speed of the classifier

based on this framework is throttled by the complexity of the architecture, which yielded a decent accuracy at the sacrifice of compute efficiency.

## Proposed Method

Leveraging the WSDM dataset, we decided to perform music preference prediction using a simple prediction model. The dataset will first go through a cleaning process to fill in empty data with a label of "-1" to indicate that the field is unavailable. Then the training dataset is split into training set and validation set with a 4:1 ratio, meaning 20% of the training data is used as the validation data. Below is the overview of the model architecture:
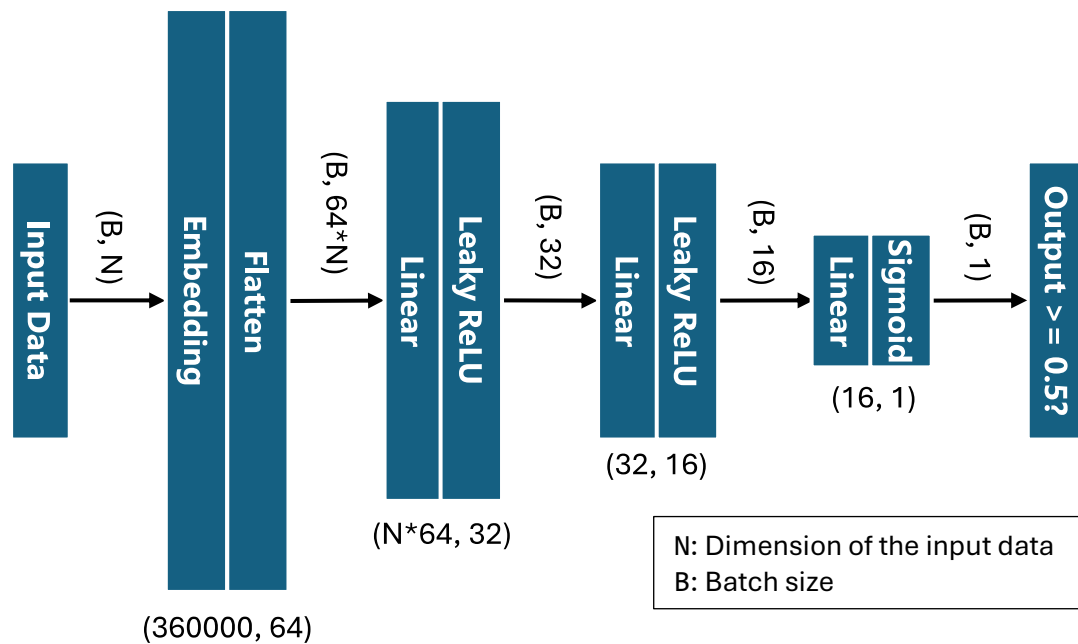


**Figure 1.** Overview of the proposed model.

After data pre-processing – further elaborated below – the input will be fed into the embedding layer and go through 3 linear layers of various neuron size and eventually output with a dimension of 1 with the sigmoid activation function. We determine that the user is likely to listen to the song again if the score exceeds 0.5, and vice versa.

*Additional model parameters* | Specified in main.py

- Optimizer: Adam
- Batch size: 800,000
- Learning rate: 0.1 (1e-1), 0.01 (1e-2), 0.001 (1e-3), 0.0001 (1e-4)
  | Learning rate scheduler: `torch.nn.ReduceLROnPlateau` (Patience: 3; Factor: 0.1)
- Number of epochs: 100 (Early stopping epochs: 10)
- Loss function: Binary Cross Entropy Loss (`torch.nn.BCELoss`)
  | Minimum improvement in validation loss to prevent early stopping: 0.01

## Results

Upon testing, the model with learning rate `1e-1` does not converge at all, achieving the minimum accuracy of around 50%. Apart from this, the model with learning rate `1e-3` continued to grow gradually beyond the 20[th] epoch, in which the other models come to an early stop. It performed similarly as the model with learning rate `1e-2` despite beating the other models that grew more sharply with half the epochs, peaking at 70.15% of accuracy.
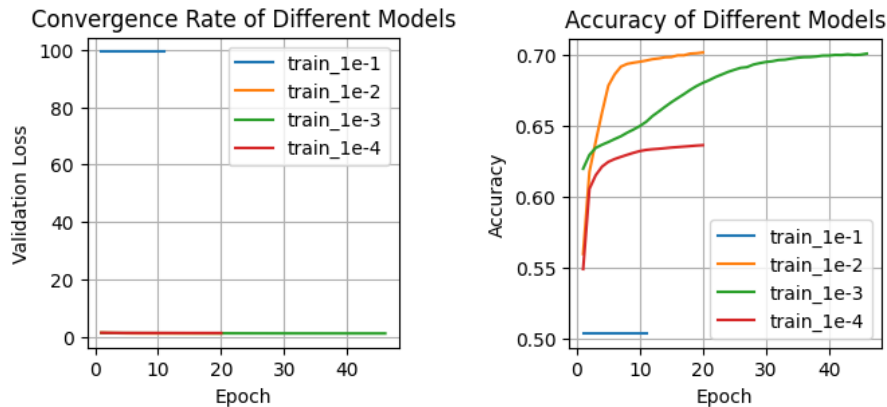


**Figure 2.** (a) Convergence rates and (b) accuracies of the model with different learning rates.

## Dataset and Code Description

src/consolidate.py

Since the WSDM dataset is split into five .csv data files, including both the training and testing data, we use this file to store functions that fuse these data into unified Pandas DataFrames, one for training and one for testing, each of which consists of: track title, Int'l Standard Recording Code (ISRC) of the track (used to record the release year of the song and to identify duplicates), user information (including listen time tracked, genre preference, etc.), genre ID, and from where user decides to play the track (i.e. which interfaces do the user go through to begin playing the song). Several attributes from the dataset are then appended to these DataFrames in the file(s) below for future use.

src/transform.py

This file houses functions that process the previously obtained DataFrames by adding higher-level attributes to the data, as well as converting values in several fields into more machine-friendly formats:

- [Line 20] We added a field from the dataset to indicate whether the artist is featured by the platform as a trending one.

- [Lines 22-29] We added fields to indicate artists with more than one role in the production of the track, specifically if s/he is also the composer and/or the lyricist, as recorded on the dataset.

- [Lines 33-43; 45-55] For tracks with multiple kinds of genre or artists, we have separated them into 3 respective fields, the first two of which label the first two items of each (genre ID and artist name), and the last labels whether more are involved.

- [Lines 57-] Finally, we encode the remaining labels in the DataFrames using `sklearn.preprocessing.LabelEncoder` to facilitate the model training.

src/net.py

This file specifies the architecture of the FFNN model used in the project.

main.py, demo.py and aws_inference_service.py

At the heart of our project are these files, responsible for model training (main.py) and deployment (the other two). The latter two, powered by `flask`, are responsible for hooking up the models on AWS for the cloud integration part of this project, elaborated in the next section.

## Cloud Deployment

Deploying a ML project using AWS EC2 involves leveraging its scalable and flexible virtual server infrastructure to meet diverse computational needs.

Deploying our machine learning project on AWS EC2 involved using the `t2.large` instance, which offers powerful CPUs ideal for training complex models. We began by launching the instance through the AWS Management Console, selecting Ubuntu as the operating system for its compatibility with ML frameworks. After setting up secure access via SSH and configuring the necessary security groups, we installed essential libraries such as PyTorch, NumPy, scikit-learn, etc. While the training process may be expedited using GPUs, given the relatively low complexity of the model, it is acceptable to be ran on a CPU-only instance.



**Figure 3.** Training process.

Throughout the model training process, we constantly monitor our machine learning models by analysing key performance metrics such as accuracy, precision and recall score using the validation dataset.



**Figure 4.** Demonstration of the Flask application running on AWS EC2 service and the batch processing function to submit the job locally.

After the training process, we have developed a simple flask-based API application to handle requests on the EC2 instance. As demonstrated in Figure 4, the API will respond to users' POST request to batch process the input data. "demo.py" includes a subset of 20 entries of the validation set just for demonstration purposes.

## Discussion

- Architecture of classification model: We have come across various machine learning architectures that can effectively extract relevant features from the dataset and perform genre classification, such as FFNN and LSTM models. We have decided to utilise the former for its simplicity and effectiveness in task accomplishment with reasonable demand for computational resources.

- AWS EC2 and ECS: Having experimented with both deployment methods, we have encountered challenges in running our model on an ECS instance owing to the limited permission of the student account. Hence, we chose to deploy our model using EC2 due to its reliability and flexibility in scaling.

## Future work

- Edge computing: By leveraging edge computing, a more lightweight version of the recommendation model can be enabled to more users using commodity devices, reducing latency and allowing the system to operate even if the user's device has limited connectivity or NPU performance. This could involve techniques like model compression, quantization, or federated learning to optimize the model for deployment on edge devices.

- Serverless model: Serverless architectures, such as those provided by cloud functions or serverless containers, can offer benefits like automatic scaling, reduced infrastructure management, etc. This could be particularly useful for handling the variable workloads and resource requirements of the music recommendation system, allowing it to scale up or down as needed without the overhead of managing a traditional server infrastructure.

## Conclusion

Having experimented with cloud-powered, ML-based music preference prediction based on the KKBOX WSDM dataset, we have gained many useful insights from the project, from model fine-tuning to hands-on experience of model deployment on AWS EC2. While we have successfully compacted the model to fit on the t2.large instance, aside from further optimizing resource usage, harnessing more powerful models can slash the training time and expedite experimentation of various architectures and model sizes.

# References

[1] KKBOX, "WSDM - KKBox's Music Recommendation Challenge," 2017. [Online] Available: https://www.kaggle.com/competitions/kkbox-music-recommendation-challenge/data

[2] W. Xing, A. Slowik and J. Peter, "Edge-cloud computing oriented large-scale online music education mechanism driven by neural networks," *J Cloud Comp,* no. 13, p. 55, 7 March 2024. doi: 10.1186/s13677-023-00555-y

# Appendix 1. Genre ID mappings

While KKBOX did not reveal the meanings of genre IDs in the WSDM dataset, the genre information embedded in these data can be decoded by comparing the genre of the track entries with the same genre ID. Below are the mappings of genre information used in the project that cover over 95% of entries in the dataset:

| ID | 139 | 359 | 437 | 444 | 451 | 458 | 465 | 726 |
|-------|------|-------|-----------|---------|-----------|-------|-------|--------|
| Genre | R&B | Soul | J-pop | K-pop | Mandopop | Pop | Rock | Rhymes |
| ID | 786 | 921 | 958 | 1209 | 1609 | 1616 | 2022 | 2122 |
| Genre | Folk | Score | Classical | Hip-hop | Electronic | House | Metal | Jazz |