

# Session-Based Cross Domain Recommendation System

Intro to Data Science, CAP5771 - Project

Santhi Daggubati

UFID: 65589133

# Milestone 1: Data Collection & Preprocessing

## Data Collection:

Three datasets were used of different domains: Movies, Music & Books from *Hugging Face* Dataset collection.

- Movies: [IMDb Dataset](#)
- Music: [Spotify Tracklist Dataset](#)
- Books: [GoodReads Dataset](#)

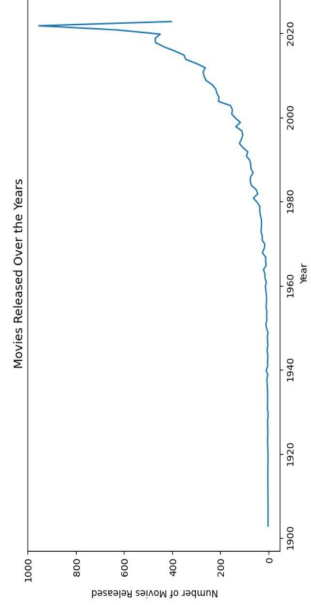
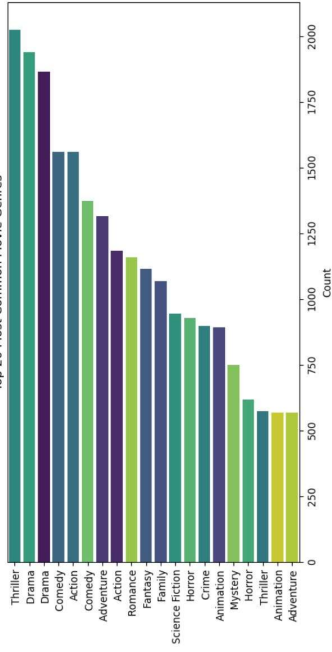
## Preprocessing:

- **Unifying Datasets:** Assigned **unique Item\_IDs** to each domain (Movies, Music, Books) by prefixing IDs with respective content type (**Movie\_**, **Music\_**, **Book\_**).
- **Renaming for Consistency:** **orig\_title** → **Title**, **genre** → **Genre**
- **Handled Missing Data**(filling and dropping), **Standardizing Columns**(consistent columns across all datasets), **Adding Item\_type column**(to specify content type: "Movie", "Music", or "Book")

# Milestone 1: Visualizations

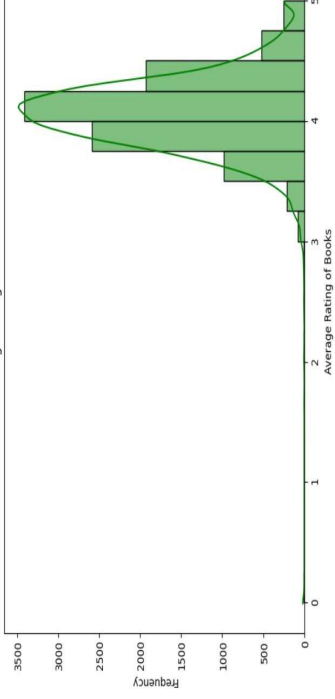
## Movies Dataset Trends:

Top 20 Most Common Movie Genres



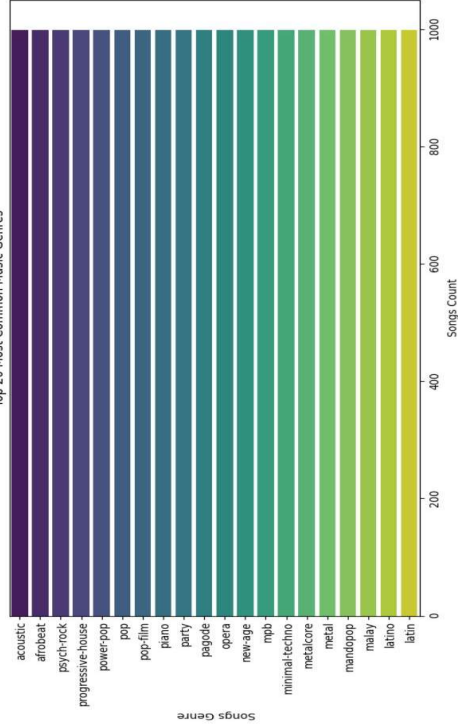
## Books Dataset Trends:

Distribution of Average Ratings of different Genre Books

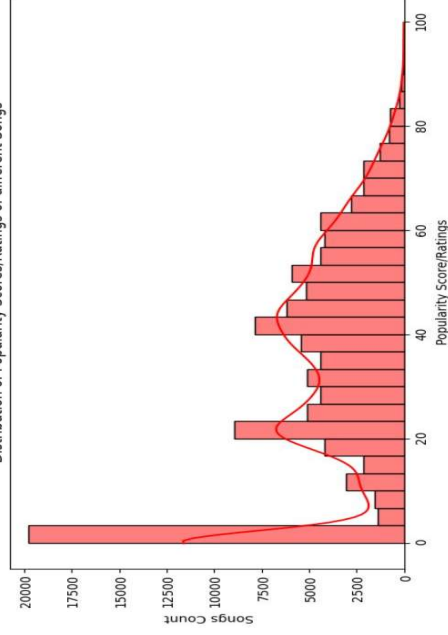


## Music Dataset Trends:

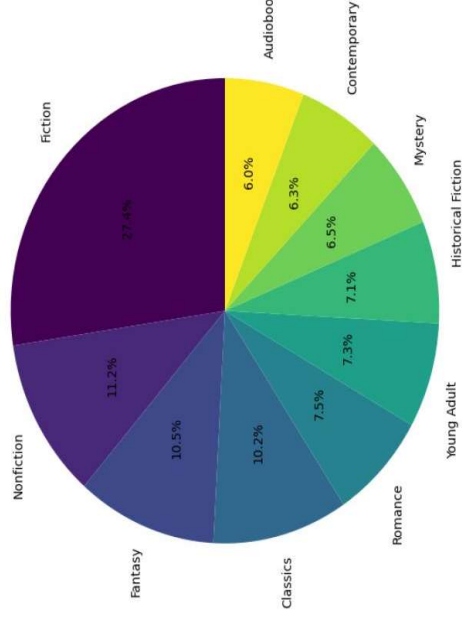
Top 20 Most Common Music Genres



Distribution of Popularity Scores/Ratings of different Songs



Top 10 Genre Distribution



# Milestone 2: Model Selection, Training & Evaluation

## Model Selection: GRU4Rec

- **Architecture:** GRU4Rec (Gated Recurrent Units)
  - **Embedding Layer:** Converts item IDs to dense vectors.
  - **GRU Layer:** Captures session-based dependencies.
  - **Fully Connected Layer:** Maps GRU output to item probabilities.
- **Reason:** Suitable for sequential recommendations based on past user interactions.

## Model Training:

- **Framework:** PyTorch
- **Training Setup:**
  - **Optimizer:** Adam (lr = 0.001)
  - **Loss:** Cross-Entropy Loss
  - **Epochs:** 100 with early stopping
  - **Batch Size:** 64
- **Evaluation:** Evaluated using **Precision@K**, **Recall@K**, and **F1 Score@K** (k - top k recommendations)

## Evaluation metrics & results:

- **Precision@10:** 0.86
- **Recall@10:** 0.89
- **F1 Score@10:** 0.87

```
# Example session (list of previously interacted Item_IDs)
example_session = ["Movie_Creed III", "Movie_Mummies", "Movie_Supercell"]

# Get top 5 recommendations
recommendations = recommend_next_items(example_session, top_k=5)
print("Recommended Items:", recommendations)
```

Recommended Items: ['Movie\_Ford v Ferrari', 'Music\_5awLjpmW05TpXCytpvCBbs', 'Music\_29RiulWABWHctRLKdqVC11', 'Music\_69Jv0CiMlppfjh9N2MFr0', 'Music\_40c70OwDazc0h2QrZhw1']

# Milestone 3: Model Integration and Deployment

## Future Work

- Tune **hyperparameters** of the model for optimal performance(better F1-Score).
- Experiment with **hybrid models** (e.g., combining collaborative filtering.)
- Develop an application using React (frontend) and Flask (backend) to deploy the recommendation system and serve real-time recommendations.
- Use Large Language Models (LLMs) to create more contextually relevant Item\_IDs for better cross-domain integration (e.g., mapping movie, music, and book data to a unified representation).