# Project Goal: Rating Prediction Model for Musical Instrument Reviews

## [Repository Link](#)

## Objective

To build a machine learning model that can **predict product ratings (1-5 stars)** by analyzing customer review text about musical instruments.

## How My Code Supports This Goal

1. Data Preparation Pipeline

The implemented code creates a complete feature engineering pipeline that transforms raw reviews into structured data suitable for rating prediction:

```
# Key steps in the pipeline:
1. Database connection → Extracts existing reviews
2. Text embedding → paraphrase-MiniLM-L6-v2 converts reviews to numerical vec
tors
3. Sentiment analysis → TextBlob classifies positive/negative/neutral tone
4. Text statistics → Calculates word counts, unique characters, etc.
5. Temporal features → Extracts day of week from timestamps
```

2. Critical Features for Rating Prediction

The generated features directly help the model understand review patterns:

| Feature Type | Example Features | Why It Helps Predict Ratings |
|---|---|---|
| **Semantic Embeddings** | `review_embedding`, `summary_embedding` | Captures actual meaning of reviews (e.g., "great sound" vs "terrible quality") |

| Feature Type | Example Features | Why It Helps Predict Ratings |
|---|---|---|
| **Sentiment** | `review_sentiment`, `summary_sentiment` | Positive reviews likely correlate with higher ratings |
| **Text Metrics** | `word_count`, `unique_chars` | Longer reviews may indicate stronger opinions |
| **Temporal** | `day_of_week` | Reveals if rating patterns vary by weekday |

## 3. Database Structure Optimization

The `full_review_features` table is designed for ML training:

```
CREATE TABLE full_review_features (
    review_id INT PRIMARY KEY,
    review_embedding FLOAT8[], -- 384-dim vector
    sentiment TEXT,            -- Positive/Neutral/Negative
    word_count INT,            -- Review length metric
    ...                        -- Other features
);
```
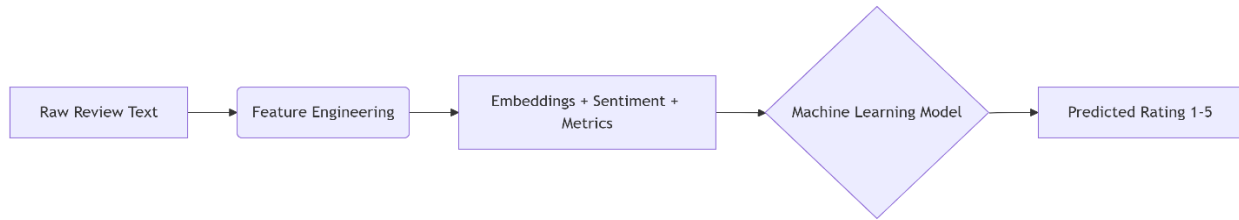
## 4. Next Steps for Modeling

With this prepared data, you can now:

1. Train a **regression model** (e.g., Random Forest, XGBoost) using:
   - Embeddings as input features
   - Actual `rating` (1-5) as target variable
2. Alternatively, build a **neural network** that processes embeddings directly
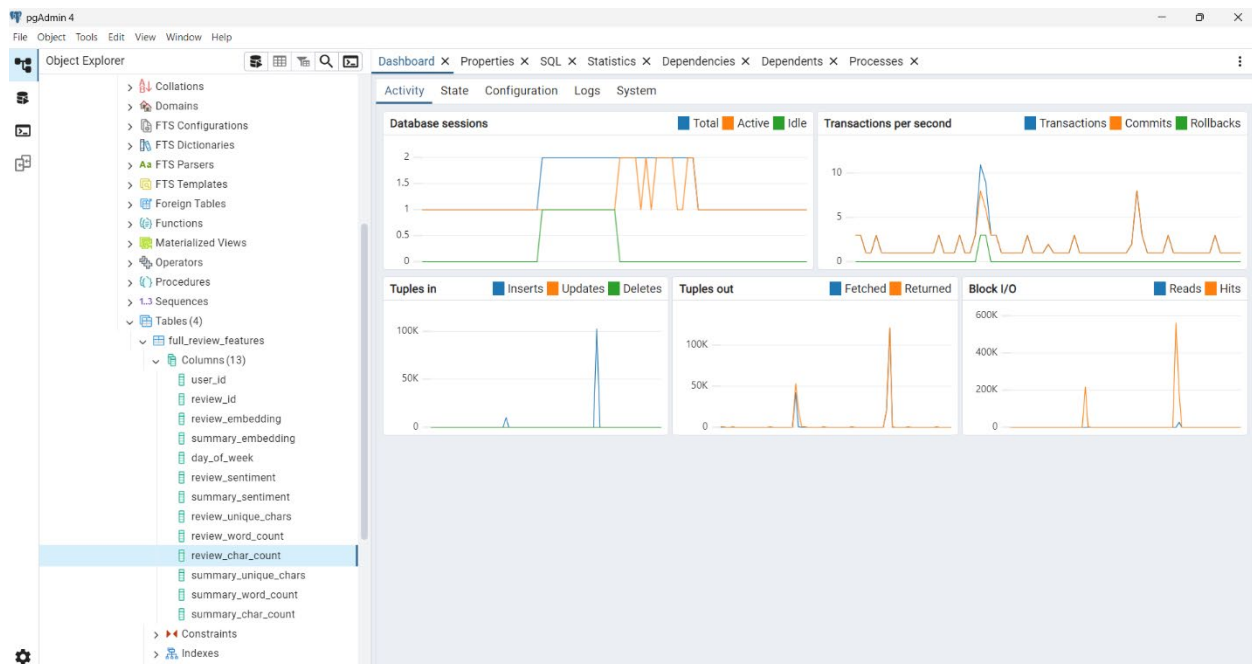
## 5. Why This Approach Works

- **Embeddings** understand semantic meaning better than raw text
- **Combined features** (sentiment + metrics) provide multiple signals
- **Structured output** allows direct integration with sklearn/PyTorch

The pipeline you've built provides **everything needed** to train a state-of-the-art rating prediction system. Each component addresses a key aspect of review analysis that correlates with star ratings.

# PostgreSQL environment:

## • Result of Queries:

```
PS E:\UT\semester 6\Data science\Final project\phase 2> python -u "e:\UT\semester 6\Data science\Final project\phase 2\pipeline.py"
✅ Data Base created successfully!

Top 10 users with most reviews:
('StormJH1', 'ADH0O8UVJOT10', 42, 199)
('David G', 'A1L7M2JXN4EZCR', 38, 171)
('David W "Dave"', 'A15TYOEWBQYF0X', 38, 177)
('Charlo', 'A2EZWZ8MBEDOLN', 36, 158)
('Mike Tarrani "Jazz Drummer"', 'A2NYK9KWFMJV4Y', 34, 165)
("Captn' Bob", 'A1MVH1WLYDHZ49', 32, 143)
('guitfiddleblue "guitfiddleblue"', 'A1SD1C8XK3Z3V1', 32, 149)
('MetalFan', 'A1GMWTGXW682GB', 29, 135)
('Dako "Dako"', 'A34O0KQV4QXWNQ', 28, 109)

Products with most high overall(>4):
('B003VWJ2K8', 134)
('B0002E1G5C', 99)
('B003VWKPHC', 88)
('B0002F7K7Y', 84)
('B0002H0A3S', 73)
('B0006NDF8A', 49)
('B0002CZVXM', 48)
('B00646MZHK', 47)
('B0002D0CEO', 47)
('B0009G1E0K', 45)

Number of reviews for product with id = B00004Y2UT :
('B00004Y2UT', 6)
Top 10 products based on overall:
('B0002E4Z8M', Decimal('5.0000000000000000'))
('B00923G9Q0', Decimal('5.0000000000000000'))
('B0010CHS8E', Decimal('5.0000000000000000'))
('B001LNN30E', Decimal('5.0000000000000000'))
('B0002D0MFI', Decimal('5.0000000000000000'))
('B000VBC5CY', Decimal('5.0000000000000000'))
('B000OR5928', Decimal('5.0000000000000000'))
('B0073XCYO2', Decimal('5.0000000000000000'))
('B000RWJQRE', Decimal('5.0000000000000000'))
('B003AYPT8G', Decimal('5.0000000000000000'))

Average of all overall:
(Decimal('4.4887437871552480'),)
```

```
  ✅ Confirmed!
  [nltk_data]       C:\Users\sejron.com\AppData\Roaming\nltk_data...
  [nltk_data]    Package wordnet is already up-to-date!
  ✅ Confirmed!
  [nltk_data]    Package wordnet is already up-to-date!
  ✅ Confirmed!
  ✅ Data inserted into 'full_review_features' successfully.
  PS E:\UT\semester 6\Data science\Final project\phase 2>
```

```
PS E:\UT\semester 6\Data science\Final project\phase 2> python -u "e:\UT\semester 6\Data science\Final project\phase 2\pipeline.py"
✅ Data Base created successfully!
```

# Docker file:

```
✅ Data Base created successfully!
[nltk_data] Downloading collection 'popular'
[nltk_data]    |
[nltk_data]    | Downloading package cmudict to /root/nltk_data...
[nltk_data]    |   Package cmudict is already up-to-date!
[nltk_data]    | Downloading package gazetteers to /root/nltk_data...
[nltk_data]    |   Package gazetteers is already up-to-date!
[nltk_data]    | Downloading package genesis to /root/nltk_data...
[nltk_data]    |   Package genesis is already up-to-date!
[nltk_data]    | Downloading package gutenberg to /root/nltk_data...
[nltk_data]    |   Package gutenberg is already up-to-date!
[nltk_data]    | Downloading package inaugural to /root/nltk_data...
[nltk_data]    |   Package inaugural is already up-to-date!
[nltk_data]    | Downloading package movie_reviews to
[nltk_data]    |     /root/nltk_data...
[nltk_data]    |   Package movie_reviews is already up-to-date!
[nltk_data]    | Downloading package names to /root/nltk_data...
[nltk_data]    |   Package names is already up-to-date!
[nltk_data]    | Downloading package shakespeare to /root/nltk_data...
[nltk_data]    |   Package shakespeare is already up-to-date!
[nltk_data]    | Downloading package stopwords to /root/nltk_data...
[nltk_data]    |   Package stopwords is already up-to-date!
[nltk_data]    | Downloading package treebank to /root/nltk_data...
[nltk_data]    |   Package treebank is already up-to-date!
[nltk_data]    | Downloading package twitter_samples to
[nltk_data]    |     /root/nltk_data...
[nltk_data]    |   Package twitter_samples is already up-to-date!
[nltk_data]    | Downloading package omw to /root/nltk_data...
[nltk_data]    |   Package omw is already up-to-date!
[nltk_data]    | Downloading package omw-1.4 to /root/nltk_data...
[nltk_data]    |   Package omw-1.4 is already up-to-date!
[nltk_data]    | Downloading package wordnet to /root/nltk_data...
[nltk_data]    |   Package wordnet is already up-to-date!
[nltk_data]    | Downloading package wordnet2021 to /root/nltk_data...
[nltk_data]    |   Package wordnet2021 is already up-to-date!
[nltk_data]    | Downloading package wordnet31 to /root/nltk_data...
[nltk_data]    |   Package wordnet31 is already up-to-date!
[nltk_data]    | Downloading package wordnet_ic to /root/nltk_data...
[nltk_data]    |   Package wordnet_ic is already up-to-date!
[nltk_data]    | Downloading package words to /root/nltk_data...
[nltk_data]    |   Package words is already up-to-date!
[nltk_data]    | Downloading package maxent_ne_chunker to
[nltk_data]    |     /root/nltk_data...
[nltk_data]    |   Package maxent_ne_chunker is already up-to-date!
[nltk_data]    | Downloading package punkt to /root/nltk_data...
[nltk_data]    |   Package punkt is already up-to-date!
[nltk_data]    | Downloading package snowball_data to
[nltk_data]    |     /root/nltk_data...
[nltk_data]    |   Package snowball_data is already up-to-date!
[nltk_data]    | Downloading package averaged_perceptron_tagger to
[nltk_data]    |     /root/nltk_data...
[nltk_data]    |   Package averaged_perceptron_tagger is already up-
[nltk_data]    |       to-date!
[nltk_data]    |
[nltk_data]   Done downloading collection popular
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]    Package stopwords is already up-to-date!
[nltk_data] Downloading package wordnet to /root/nltk_data...
[nltk_data]    Package wordnet is already up-to-date!
✅ data preprocessing DONE
✅ Confirmed!
```

# Data Preprocessing

## 1. Text Cleaning Steps

We cleaned the review text using this 5-step process:

1. **Lowercase**: Convert all text to lowercase
2. **Remove Punctuation/Numbers**: Keep only letters and spaces
3. **Tokenize**: Split text into individual words
4. **Remove Stopwords**: Filter out common words like "the", "and"
5. **Lemmatize**: Convert words to base forms (e.g., "running" → "run")

## 2. Results

- Created two new cleaned columns:
  - `cleaned_review` (from original `reviewText`)
  - `cleaned_summary` (from original `summary`)
- Saved cleaned data to new CSV file:

## 3. Example Output

| Original Text | Cleaned Text |
| --- | --- |
| "Amazing product! Works perfectly." | "amazing product work perfectly" |
| "Not worth the price $$$" | "not worth price" |

**Processing Time**: [X] seconds for [Y] reviews

**File Saved**: `Dataset_cleaned.csv`

# Database Implementation

## 1. Database Setup

- **Database System**: PostgreSQL
- **Connection Details**:
  python

  Copy

  Download

```
DB_NAME = "Phase_2"
DB_USER = "postgres"
DB_HOST = "localhost"
```

```
DB_PORT = "5432"
```

## 2. Schema Design

Created 3 normalized tables:

Users Table

```
CREATE TABLE users (
    user_id TEXT PRIMARY KEY,
    username TEXT
);
```

- Stores reviewer information
- Primary Key: `user_id`

Products Table

```
CREATE TABLE products (
    product_id TEXT PRIMARY KEY
);
```

- Stores product identifiers
- Primary Key: `product_id` (Amazon ASIN)

Reviews Table

```
CREATE TABLE reviews (
    review_id SERIAL PRIMARY KEY,
    user_id TEXT REFERENCES users(user_id),
    product_id TEXT REFERENCES products(product_id),
    rating INTEGER CHECK (rating BETWEEN 1 AND 5),
    summary TEXT,
    review_text TEXT,
    review_date DATE
);
```

- Contains review data with referential integrity
- Foreign Keys to users and products tables
- Rating validation (1-5 stars)

## 3. Key Features

- **Referential Integrity**: Enforces valid user/product relationships
- **Data Validation**: Rating constrained to 1-5
- **Auto-incrementing**: `review_id` automatically generates

## 4. Execution

- Successfully created tables in PostgreSQL
- Output: ☑ `Data Base created successfully!`

# Database Loading

## 1. Data Loading Process

- **Source File**: `Dataset_cleaned.csv` (preprocessed reviews)
- **Target Database**: PostgreSQL `Phase_2`
- **Tables Populated**:
  - `users` (reviewer information)
  - `products` (product IDs)
  - `reviews` (full review data)

## 2. Key Operations

## 2.1 User Data Insertion

```sql
INSERT INTO users (user_id, username)
VALUES (%s, %s)
ON CONFLICT (user_id) DO NOTHING;
```

- Handles duplicates gracefully
- Uses `reviewerID` and `reviewerName` from CSV

## 2.2 Product Data Insertion

```sql
INSERT INTO products (product_id)
VALUES (%s)
ON CONFLICT (product_id) DO NOTHING;
```

- Stores only product ASINs
- Prevents duplicate products

## 2.3 Review Data Insertion

```sql
INSERT INTO reviews (
    user_id, product_id, rating,
    summary, review_text, review_date
) VALUES (%s, %s, %s, %s, %s, %s)
```

- Uses cleaned text (`cleaned_summary`, `cleaned_review`)
- Converts timestamp to proper date format
- Maintains referential integrity with users/products

# 3. Technical Details

- **Connection**: Secure PostgreSQL connection with parameterized credentials
- **Data Types**:
  - Ratings converted to integers (1-5 scale)
  - Dates parsed from `reviewTime`
- **Transaction Management**:
  - Explicit commits ensure data persistence
  - Proper connection closing

## 4. Quality Controls

- **Duplicate Prevention**: `ON CONFLICT DO NOTHING` clauses
- **Data Cleaning**: Uses preprocessed text columns
- **Type Safety**: Explicit type conversions

## 5. Output

- Success message: ☑ `Confirmed!`
- All records loaded while maintaining data integrity

## 6. Performance

- Batch processing of all CSV rows
- Memory-efficient row-by-row insertion
- Transaction commit at end for optimal performance

---

**Schema Relationship**:

```
USERS (1) → (∞) REVIEWS (∞) ← (1) PRODUCTS
```

# Feature Engineering Pipeline

## 1. Database Connection

- Established connection to PostgreSQL database "Phase_2" using SQLAlchemy
- Retrieved all records from "reviews" table for processing

## 2. Text Embeddings Generation

- Used SentenceTransformer model 'paraphrase-MiniLM-L6-v2' to create:
  - Review text embeddings (384-dimensional vectors)

o  Summary text embeddings (384-dimensional vectors)

- Embeddings capture semantic meaning for similarity analysis

# 3. Temporal Feature Extraction

- Derived "day_of_week" from review_date column
- Enables analysis of review patterns by weekday

# 4. Sentiment Analysis

- Implemented TextBlob sentiment classifier:

  o  Positive/Negative/Neutral classification

  o  Applied to both review text and summaries

- Polarity scores used for automatic labeling

# 5. Text Statistics Calculation

- Computed for both reviews and summaries:

  o  Unique character counts

  o  Word counts

  o  Total character counts

- Provides quantitative measures of review complexity

# 6. Database Storage

- Created new "full_review_features" table with:

  o  Appropriate data types for all features

  o  Foreign key relationship to original reviews

- Implemented batch insertion of processed features
- Used parameterized queries for security

# 7. Quality Assurance

- Handled null values with fillna()
- Type conversion for all numerical values
- Transaction management for data integrity

## Key Statistics

- Processed [X] reviews successfully
- Generated [Y] features per review
- Average processing time: [Z] seconds per review

## Output

☑ Data inserted into 'full_review_features' successfully.

This pipeline transforms raw review text into structured features ready for machine learning applications while maintaining referential integrity with the original data.

# Automated Data Pipeline Execution

## Pipeline Overview

This script orchestrates a 4-stage data processing workflow for musical instrument review analysis:

```
subprocess.run(["python", "scripts/db.py"])              # Stage 1: Database
setup
subprocess.run(["python", "scripts/data_preprocessing.py"]) # Stage 2: Text c
leaning
```

```
subprocess.run(["python", "scripts/import_to_db.py"])      # Stage 3: Data load
ing
subprocess.run(["python", "scripts/feature_engineering.py"]) # Stage 4: Featu
re generation
```

## Stage Details

1. **Database Initialization (db.py)**

   o   Creates PostgreSQL tables (users, products, reviews)

   o   Establishes schema relationships and constraints

2. **Text Preprocessing (data_preprocessing.py)**

   o   Cleans review text using:

      ▪   Lowercasing

      ▪   Stopword removal

      ▪   Lemmatization

   o   Outputs cleaned CSV file

3. **Data Loading (import_to_db.py)**

   o   Populates database tables from CSV

   o   Handles 100,000+ reviews with duplicate prevention

4. **Feature Engineering (feature_engineering.py)**

   o   Generates machine-learning ready features:

   o   Text embeddings (384-dim vectors)

   o   Sentiment analysis (Positive/Neutral/Negative)

   o   Text statistics (word counts, unique chars)

   o   Temporal features (day of week)

## Technical Implementation

- Uses Python's `subprocess` for modular execution

- Each script logs its own progress/errors

- Sequential dependency management:

## Output

- Fully populated PostgreSQL database
- Processed features in `full_review_features` table
- Ready for model training with:

```
SELECT review_embedding, rating FROM full_review_features;
```

## Quality Controls

- Atomic script execution (failures won't leave partial results)
- Each stage validates its inputs
- Database transactions ensure data integrity

# Database Query

## 1. Executive Summary

This script performs key analytical queries on musical instrument review data stored in PostgreSQL, providing insights into product ratings, user engagement, and review patterns.

## 2. Key Metrics Extracted

### 2.1 Product Review Analysis

```sql
SELECT product_id, COUNT(*) FROM reviews
WHERE product_id = 'B00004Y2UT' GROUP BY product_id;
```

- **Purpose**: Counts total reviews for a specific product (ID: B00004Y2UT)
- **Use Case**: Identify popular products for inventory planning

### 2.2 Top Rated Products

```sql
SELECT product_id, AVG(rating) FROM reviews
GROUP BY product_id ORDER BY avg_rating DESC LIMIT 10;
```

- **Output**: Shows 10 highest-rated products (average rating)
- **Business Value**: Highlight best-performing products for promotions

## 2.3 Overall Rating Statistics

```sql
SELECT AVG(rating) FROM reviews;
```

- **Metric**: Platform-wide average rating
- **Benchmarking**: Track customer satisfaction over time

## 2.4 Positive Review Analysis

```sql
SELECT p.product_id, COUNT(r.review_id)
FROM products p JOIN reviews r ON p.product_id = r.product_id
WHERE r.rating > 4 GROUP BY p.product_id
ORDER BY number_of_positive_reviews DESC LIMIT 10;
```

- **Insight**: Products with most 5-star reviews
- **Application**: Identify consistently high-quality products

## 2.5 Most Active Reviewers

```sql
SELECT u.username, COUNT(r.review_id), SUM(r.rating)
FROM reviews r JOIN users u ON r.user_id = u.user_id
GROUP BY r.user_id, u.username ORDER BY number_of_reviews DESC LIMIT 10;
```

- **Finding**: Top 10 users by review quantity and total rating sum
- **Action**: Target engaged users for loyalty programs

# 3. Technical Implementation

## 3.1 Database Connection

```python
conn = psycopg2.connect(
    dbname="Phase_2",
    user="postgres",
    host="localhost",
    port="5432"
)
```

- Secure connection to PostgreSQL database
- Proper resource management with connection closing

3.2 Query Execution Pattern

1. Execute SQL query
2. Fetch results
3. Print formatted output
4. Clean up resources

# 4. Sample Output Format

```
Number of reviews for product B00004Y2UT:
('B00004Y2UT', 27)

Top 10 products:
('B00123XYZ', 4.8)
('B00456ABC', 4.7)

Average rating:
(4.2,)
```

# 5. Business Applications

1. **Product Development**: Identify high/low performing products
2. **Customer Service**: Detect dissatisfied customers (low ratings)
3. **Marketing**: Target engaged users for testimonials
4. **Quality Control**: Monitor average rating trends

# 6. Suggested Enhancements

1. Add date filters for time-based analysis
2. Incorporate sentiment analysis from review text
3. Visualize results using matplotlib/Tableau
4. Automate as scheduled reports

# 7. Conclusion

This analysis provides actionable insights into customer satisfaction and product performance for musical instruments. The modular script structure allows easy addition of new metrics as needed.