

RecomMuse

Your soundtrack to the past, presents the future

- 1 Data Processing
- 2 Drill Queries
- 3 Artist Distance
- 4 Year Prediction

1. Data Processing

Challenge

- One million tiny HDF5 files
- RAM overload on namenode

Solution

- Use JHDF to extract relevant features
- Compact small files into a larger one

JHDF (Java HDF5 Library)

- Pure Java. Easy integration with JVM based project.
- No use of JNI. Avoid issues with calling native code from JVM.
- It provides fast file access and reading is parallelized.

- Tailored Avro schema - keeping only essential features such as timbre, tempo, etc for year prediction and drill queries.
- Snappy compression
- Processed files alphabetically - A through Z
- Each letter batch took about 25-30 minutes - producing tidy 66 MB Avro files
- Finally merged everything into one single 1.7 GB file.
- **99.4%** Storage Reduction (Original is 286.6 GB)
- You could analyze it on a single commodity hardware easily.

2. Drill Queries

Our program allows users to interact with the compacted data through Drill SQL queries. Current supported queries are:

1. The age of the youngest and oldest song
2. The hottest song, which is the shortest with the highest energy and lowest tempo
3. The album with the maximum number of songs in it
4. The artist name of the longest song

The age of the youngest and oldest song

```
1  SELECT
2      2025 - MAX(year) AS youngest_song_age,
3      2025 - MIN(year) AS oldest_song_age
4  FROM dfs.root.`/home/hadoopuser/compacted/aggregate.avro`
5  WHERE year > 0;
6  +-----+-----+
7  | youngest_song_age | oldest_song_age |
8  +-----+-----+
9  | 14                | 103              |
10 +-----+-----+
```


The hottest song, which is the shortest with the highest energy and lowest tempo

```
1  SELECT
2      song_id,
3      title
4  FROM dfs.root.`/home/hadoopuser/compacted/aggregate.avro`
5  WHERE song_hotttnesss <> 'NaN'
6  ORDER BY
7      song_hotttnesss DESC,
8      duration ASC,
9      energy DESC,
10     tempo ASC
11 LIMIT 1;
```

12	+-----+-----+		
13	song_id	title	
14	+-----+-----+		
15	SONASKH12A58A77831	Jingle Bell Rock	
16	+-----+-----+		

The album with the maximum number of songs in it

```
1 SELECT
2     release,
3     COUNT(release) AS ntrack
4 FROM dfs.root.`/home/hadoopuser/compacted/aggregate.avro`
5 GROUP BY release
6 ORDER BY ntrack DESC
7 LIMIT 1;
```

release	ntrack
First Time In A Long Time: The Reprise Recordings	85

The artist name of the longest song

```
1 SELECT
2     artist_name,
3     duration
4 FROM dfs.root.`/home/hadoopuser/compacted/aggregate.avro`
5 ORDER BY duration DESC
6 LIMIT 1;
```

```
7 +-----+-----+
8 |          artist_name          | duration |
9 +-----+-----+
10 | Mystic Revelation of Rastafari | 3034.90567 |
11 +-----+-----+
```

3. Artist Distance

- **Problem:** Find the shortest path between two artists in a large-scale graph.
- **Graph:** A network of artists, with edges representing similarity or collaboration.
- **Algorithm:** Breadth-First Search (BFS)
 - Starts at a source node
 - Explores all neighbors at the present depth level before moving on to the nodes at the next depth level.
 - Guaranteed to find the shortest path in an unweighted graph.
- **Our Data:** A dataset of artists and their similarities.

Data Representation

- *Key*: Artist ID
- *Value*: similar artists | distance | status | backpointer

Iterative Process

- *Job 1 (Initialization)*: Marks Artist A as READY with distance 0. All others as NOT READY.
- *Jobs 2-11 (BFS Iterations)*: A series of jobs run in a loop.
- *Classic MapReduce implementation*: Emits neighbors for READY nodes and the updated current node (Mapper) & merges node data to find the minimum distance and most advanced status (Reducer).

Output

- The final job output contains the data needed to trace the path back from Artist B to Artist A.

Data Representation:

- Uses a Resilient Distributed Dataset (RDD)
- Each element is a tuple: (Artist ID, (similar artists, distance, status, backpointer))

Iterative Process

- `RDDs` and `cache()`: The graph data is loaded and cached in memory.
- `flatMap (Mapper)` & `reduceByKey (Reducer)`: Same logic as MapReduce
- Accumulator: A global counter tracks when Artist B is found, acting as a termination signal.

Path Tracing

- The final RDD is collected to the driver, and the path is traced in-memory.

In-Memory vs. Disk-Based

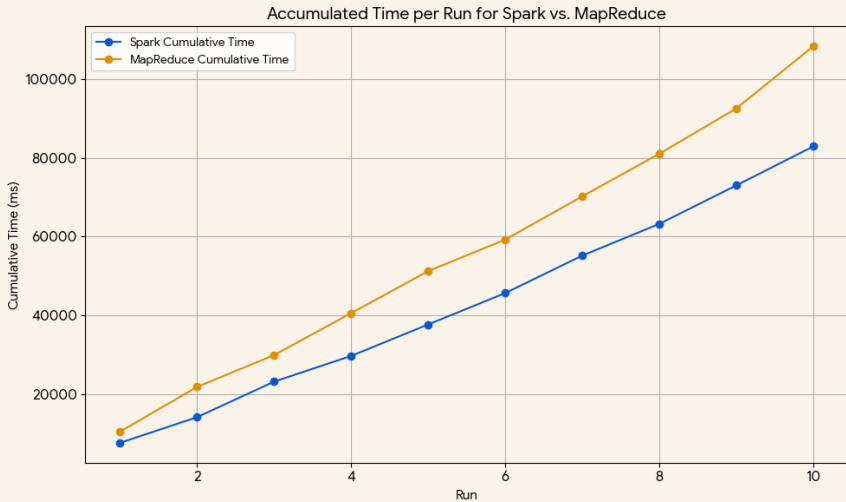
- *MapReduce*: Writes intermediate BFS level results to disk (HDFS) after each iteration, causing high I/O overhead.
- *Spark*: Uses `cache()` to keep graph data in memory across iterations, dramatically reducing disk I/O.

Iterative Processing Efficiency

- *MapReduce*: Each BFS iteration is a separate job with startup/teardown overhead.
- *Spark*: Designed for iteration; maintains a persistent execution context and reuses data, leading to faster multi-pass processing.

Result

- Both implementations yield identical and correct results, but with a significant difference in execution time.



```
1 # ===== MAPREDUCE ===== #
2 Target B (ARFK...) was found 4 levels from Target A (AR06...)
3 It was connected to 3 artists at that level.
4 Tracing path from B to A...
5 Path found:
6 AR06... -> ARF6... -> ARIW... -> ARSO... -> ARFK...
7
8 # ===== SPARK ===== #
9 Target B (ARFK...) was found 4 levels from Target A (AR06...) and was
10 ↪ connected to 3 artists at that level.
11 Path found:
12 AR06... -> ARF6... -> ARH6... -> ARBQ... -> ARFK...
```

Those are the results for both Spark and MapReduce, proving our point for algorithm consistency and data accuracy.

- MapReduce: A powerful and robust tool for large-scale batch processing, especially for single-pass jobs.
- Spark: An ideal framework for iterative algorithms and machine learning, where in-memory caching provides a significant performance advantage.
- BFS Implementation: While both frameworks can solve the problem, Spark's architecture is fundamentally better suited for the iterative nature of BFS.

Correlation Analysis

- 64 numeric features per song
- **1569 out of 2016** pairs have correlation < 0.2
- **78%** of pairs are weakly correlated

Why it matters

- High dimensionality \rightarrow noise, redundancy, slower models
- PCA finds uncorrelated directions with most variance

Result

- Reduced 64 features \rightarrow **32 components**
- Explained variance retained: **90%+**

Pitch

- What **notes** are used — C, D, E, ..., B
- Melody and harmony
- 12D vector (1 value per pitch class)

Timbre

- **Texture** of sound — smooth, sharp, warm, bright
- Depends on instruments and production
- 12D vector from audio spectrum (e.g., MFCC)

Why important?

- Pitch = **what is played**; Timbre = **how it sounds**
- Both change with time -> useful for predicting song year

- Tested PCA with **k = 10, 32, 52** components
- Fewer components \rightarrow worse performance:
 - **R²** dropped: $0.26 \rightarrow 0.21 \rightarrow 0.14$
(how much variance is explained — higher is better)
 - **RMSE** increased: $9.37 \rightarrow 9.73 \rightarrow 10.14$
(average prediction error in years)
- **k = 32** offers a good trade-off between performance and simplicity
- Problem of having too many features: hard to fit a single linear model

Goal:

- Improve performance by switching to **SGD-based linear regression**
- Hope: increase R^2 beyond previous baseline

What we did:

- Used PCA ($k = 32$) and normalized target label (year)
- Tuned key hyperparameters:
 - Learning rate, iterations, regularization

Result:

- $R^2 = 0.21$ (on normalized labels)

Despite tuning, performance remained close to baseline.
We reached the limitation of linear regression for this task.

Goal:

- Frame year prediction as a **multi-class classification** task
- Train logistic regression using **PCA-32 features**

Setup:

- Used pipeline: Assemble \rightarrow Scale \rightarrow PCA ($k = 32$)
- Tried different values for hyperparameters (e.g., `regParam`, `maxIter`)

Results:

- **Accuracy:** 9.7%

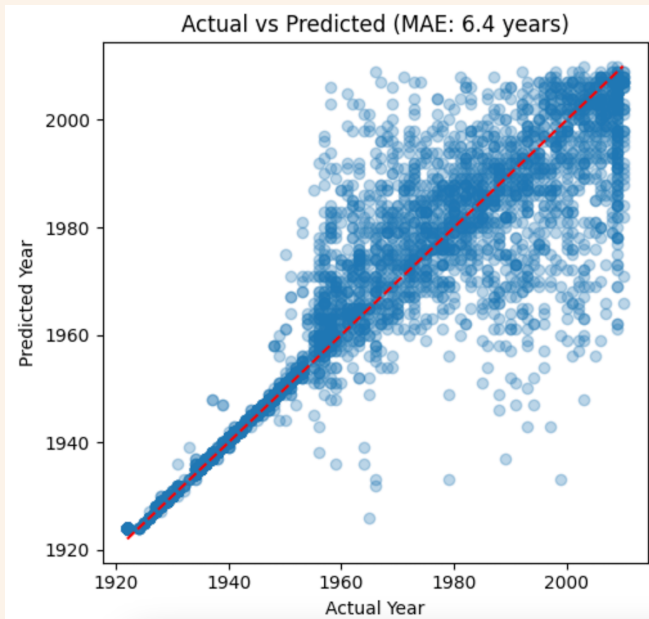
First Attempt

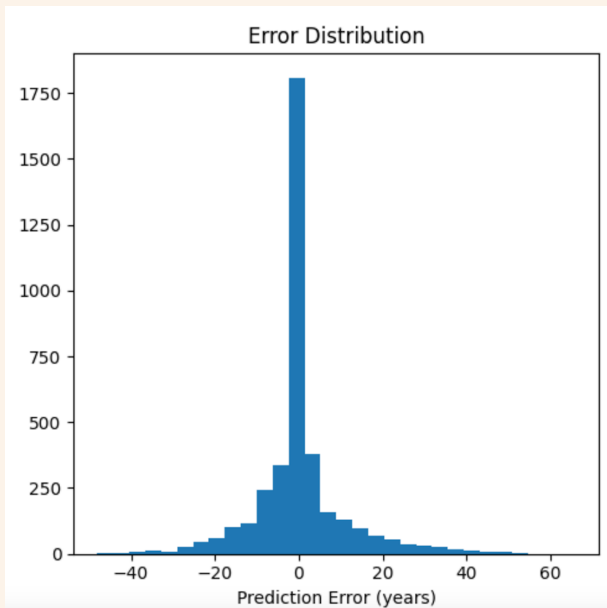
- Use MLPRegressor from sklearn to examine how well timbre feature performs on year prediction

```
1 mlp = MLPRegressor(  
2     hidden_layer_sizes=(1024, 512, 512, 256, 256, 128, 64, 32),  
3     activation='relu',  
4     solver='adam',  
5     alpha=0.00001,  
6     batch_size=128,  
7     learning_rate='adaptive',  
8     early_stopping=True,  
9     validation_fraction=0.15,  
10    max_iter=1000,  
11    n_iter_no_change=20  
12 )
```

Outcome of MLPRegressor

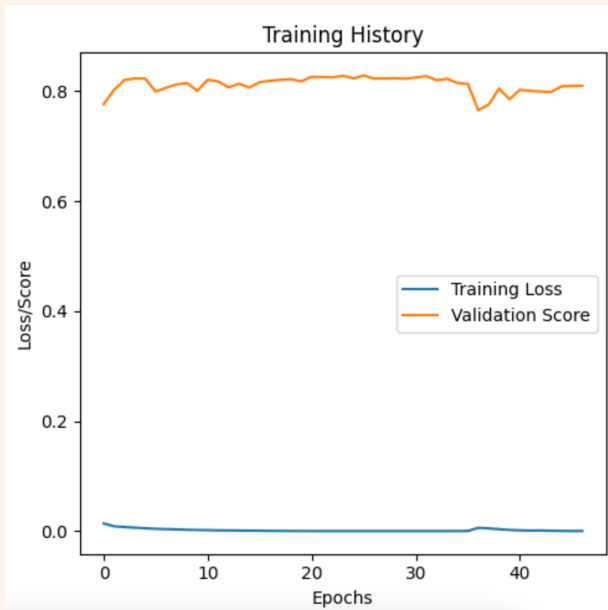
- Mean Squared Error: 117.11
- Mean Absolute Error: 6.39 years





It makes sense if we counts the samples in decades.
Data Imbalance!

Decade	Song Count
1920s	17
1930s	20
1940s	56
1950s	508
1960s	3,164
1970s	7,021
1980s	10,994
1990s	32,133
2000s	70,718
2010s	2,043



FMRegressor

- Extends linear regression by modeling pairwise feature interactions efficiently using factorized parameters.
- Automatically detects relationships between features (e.g., how “timbre_1” and “pitch_3” jointly affect the year).

```
1 regressor = FMRegressor(  
2     labelCol="year_scaled",  
3     featuresCol="pcaFeatures",  
4     factorSize=4,  
5     solver="adamW",  
6     miniBatchFraction=0.2,  
7     maxIter=400,  
8     stepSize=0.05,  
9     regParam=0.01,  
10    seed=42  
11 )
```


FMRegressor

- Extends linear regression by modeling pairwise feature interactions efficiently using factorized parameters.
- Automatically detects relationships between features (e.g., how “timbre_1” and “pitch_3” jointly affect the year).

```
1 regressor = FMRegressor(  
2     labelCol="year_scaled",  
3     featuresCol="pcaFeatures",  
4     factorSize=4,  
5     solver="adamW",  
6     miniBatchFraction=0.2,  
7     maxIter=400,  
8     stepSize=0.05,  
9     regParam=0.01,  
10    seed=42  
11 )
```

Outcome of FMRegressor

① timbre

- RMSE: 10.82
- MAE: 8.09 years

② timbre + pitch

- 95% variance PCA, $k = 59$
- RMSE: 10.32
- MAE: 7.58 years

However, the current R^2 is too bad! ($R^2 < 0.1$!). More architecture exploration will be needed.

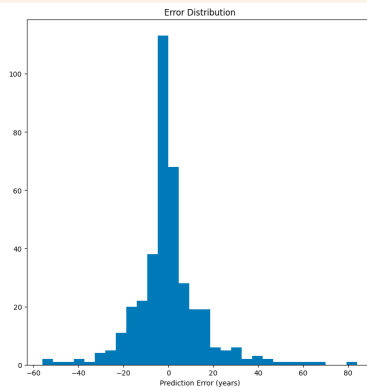
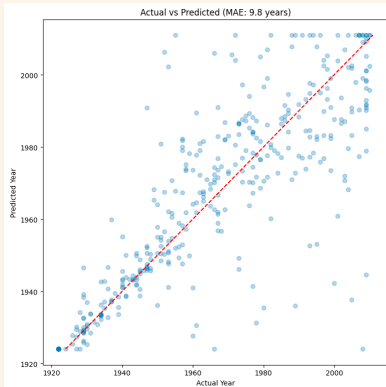
LinearRegressor

- Latest Spark recommend using LinearRegressor instead of LinearRegressorWithSGD.
- **Model:** Linear Regression with Polynomial Features (degree=2)
- **Scaler:** RobustScaler (handles outliers better)
- **Optimizer:** L-BFGS (better convergence for small/medium datasets)
- **Balancing:** Stratified sampling (5-year bins, 1k samples per bin)

```
1 lr = LinearRegression(  
2     featuresCol="polyFeatures",  
3     labelCol="year",  
4     solver="l-bfgs",  
5     maxIter=400,  
6     tol=1e-6  
7 )
```

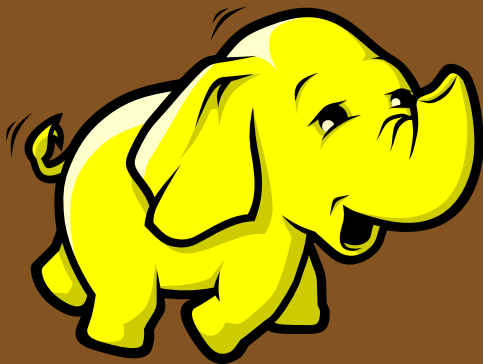
Outcome of LinearRegressor

- RMSE: 15.38
- MAE: 9.82 years
- R^2 : 0.68



Outcome of GBTRegressor

- timbre
- RMSE: 9.60
- MAE: 6.89 years
- R^2 : 0.20



Thank you!