# Meta Musical Memes

Megan Jasek, James King, Sean Underwood
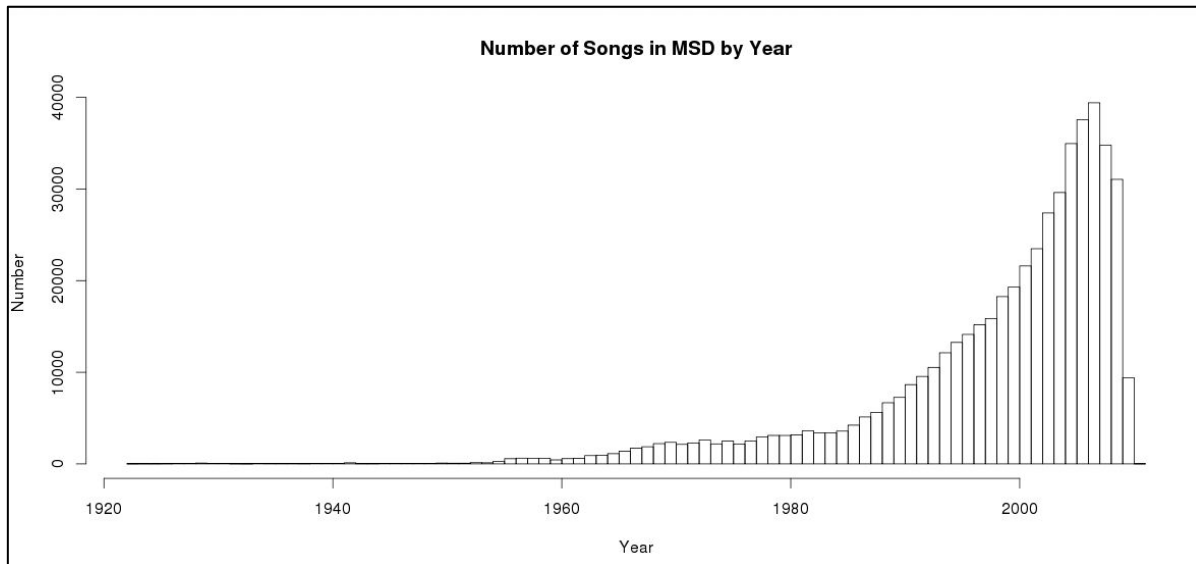
# Problem Definition

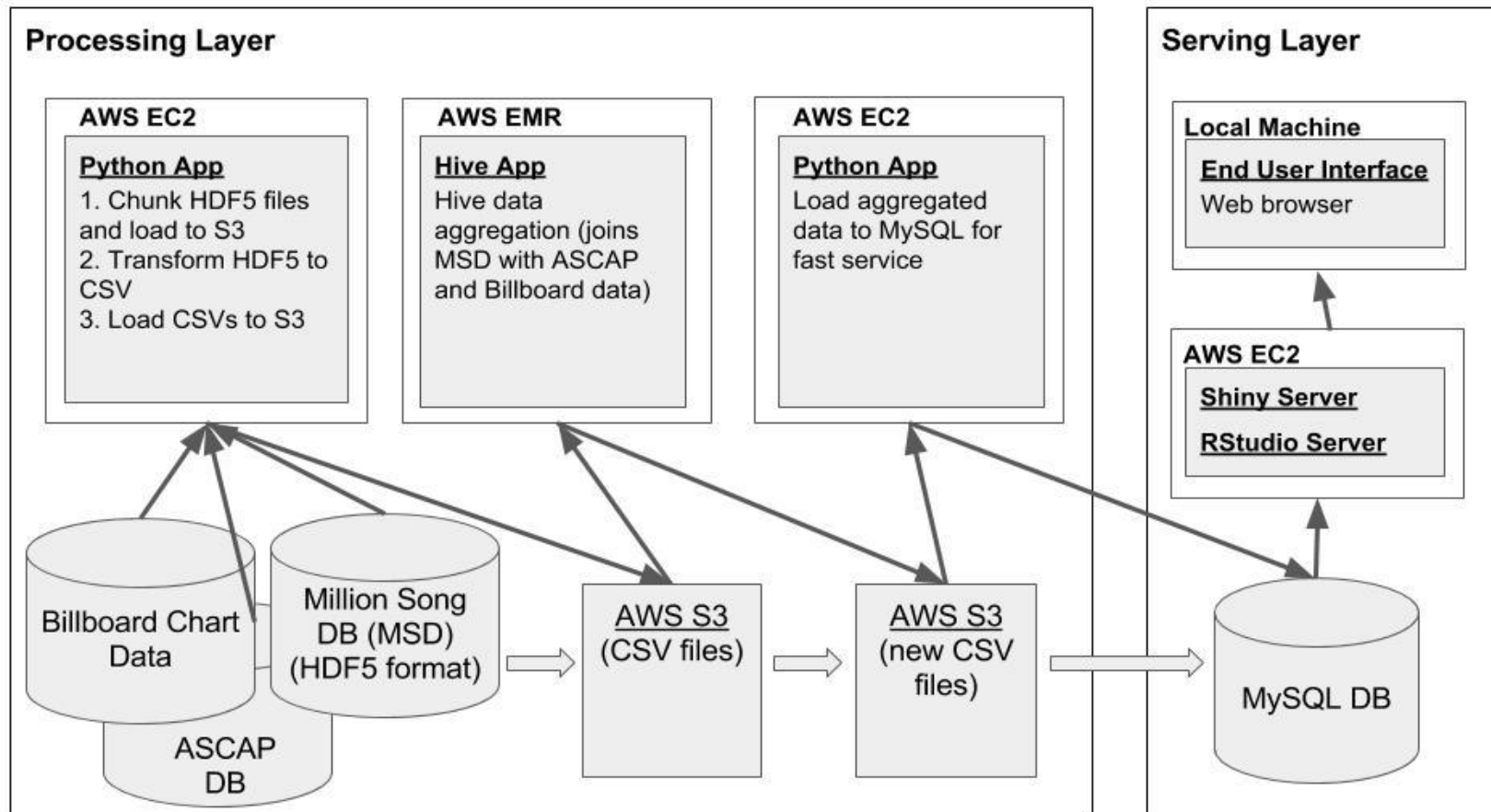Is there a basic formula that helps predict a song's commercial success?

# Data

- **Million Song Dataset**
  - Song credits
  - Music theory info (key, tempo, etc.)
  - Audio features
  - Popularity estimates ("Hotttnesss")
- **Billboard charts**
  - Performer credits
  - Peak position
  - Peak year
- **ASCAP records**
  - Writer credits

**Number of Songs in MSD by Year**

# Data - Acquisition and Organization

- Million Song Dataset:
  - Provided by Columbia University's LabROSA
  - Publicly available on Amazon AWS as an EBS Snapshot
  - File format: HDF5 (one file per song)
  - 500 GB
- Billboard Data
  - Fan-maintained
  - Excel format
  - 29 MB
- ASCAP Data
  - Provided by ASCAP
  - CSV format
  - 1.6 GB

# Architecture



**Processing Layer**

**AWS EC2**

**Python App**
1. Chunk HDF5 files and load to S3
2. Transform HDF5 to CSV
3. Load CSVs to S3

**AWS EMR**

**Hive App**
Hive data aggregation (joins MSD with ASCAP and Billboard data)

**AWS EC2**

**Python App**
Load aggregated data to MySQL for fast service

**Serving Layer**

**Local Machine**
**End User Interface**
Web browser

**AWS EC2**
**Shiny Server**
**RStudio Server**

Billboard Chart Data

Million Song DB (MSD) (HDF5 format)

ASCAP DB

AWS S3 (CSV files)

AWS S3 (new CSV files)

MySQL DB

# Processing Layer

- ETL:  Million Song Dataset (500 GB)
  - HDF5 - Does not interact well with HDFS.
    - HDF5 is optimized for high performance when being used from a single process that does not need to load the whole file into memory at once.
    - When a file gets split across a block boundary, there's no good way for workers to re-assemble the data.
    - h5py python library
  - Extraction - multiprocessing python library - run code in parallel
    - xLarge AWS instance with 8 VCPUs
    - Still took 12 hours to extract data
- ETL:  Billboard Data (30 MB)
  - MS Excel file
  - Extracted 4 attributes (song, artist, peak position, peak year) using python
- ETL:  ASCAP Data (1.5 GB)
  - CSV file in long format, unusual name formatting:  "last middle first"
  - Extracted 3 attributes (song, artist, writer)
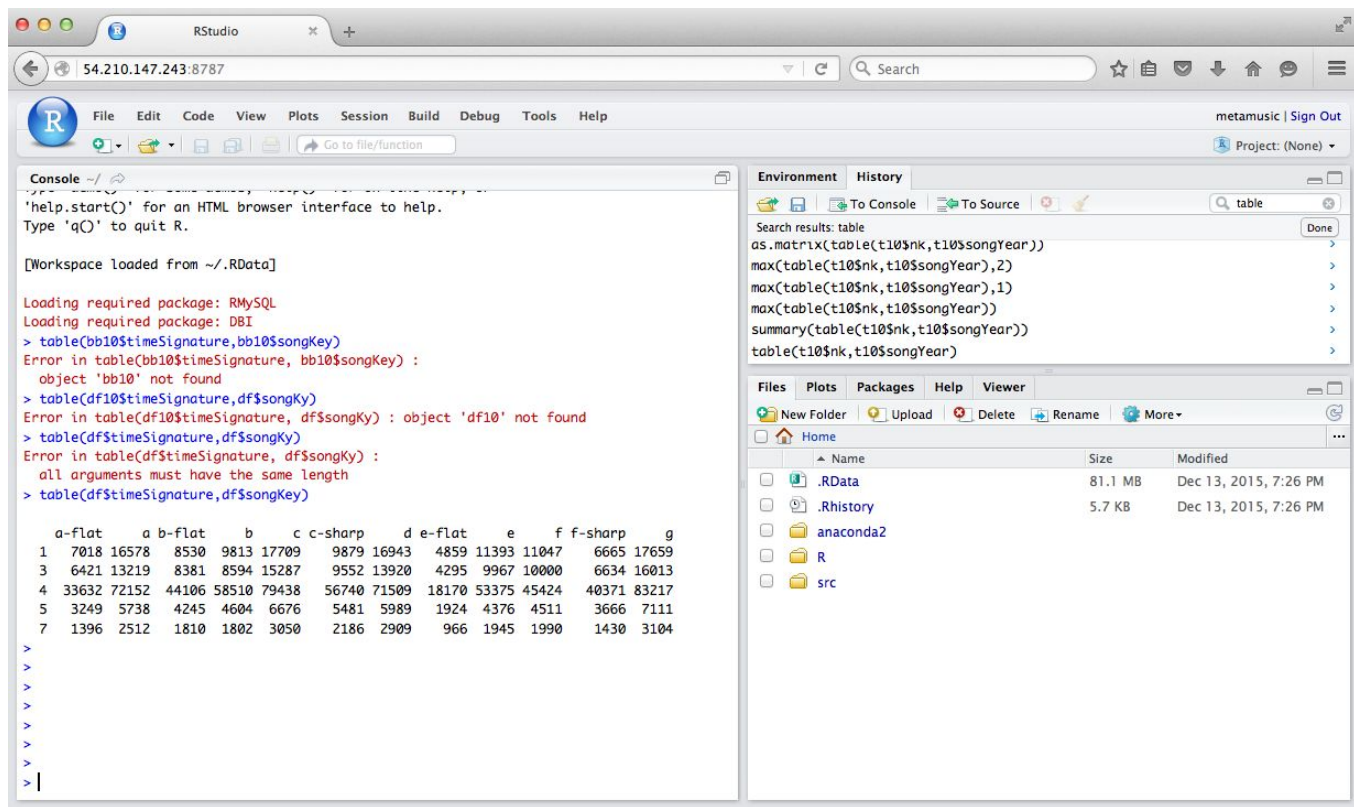
# Processing Layer (cont.)

- Data Merge
  - All data converted to CSV format
  - Uploaded to Amazon S3
  - Amazon Elastic Map Reduce (EMR)
    - Hive job to clean and merge datasets
      - Massaged MSD 0 values
      - Converted numeric encodings for key and mode to text
    - Output uploaded to Amazon S3
    - 3 m3.xlarge instances - total time 5 minutes
  - Hive trickiness
    - Hive does not proactively validate data
    - It's very important to manually validate
- ASCAP: There's no escape from incomplete data
  - We failed to find a reliable way to bridge the lack of performer credits in the ASCAP data.
- Produced 2 datasets
  - One with ASCAP and one without

# Serving Layer

- Merged data in MySQL DB
- Interactive Output
  - RStudio Server
    - Running on EC2 instance
    - Exposed full R ecosystem in a web browser
  - Shiny
    - Toolkit for writing web applications in R
    - Create interactive UI widgets hooked up to R commands
  - Shiny Server
    - Running on EC2 instance
    - Exposes Shiny applications in a web browser

# RStudio Server

# RStudio Results:  Songs are Getting Louder…..What?



Median Loudness by Year

# RStudio Results: Songs. Aren't. Getting. Any. Faster.



Median Tempo by Year

# RStudio Results: Songs are Getting Lonnnnnnger...



Median Song Duration by Year

# Analysis: Shiny Example 1:

## Song Statistics by Artist

# Analysis: Shiny Example 2:

## Song Hotness Regression



52.91.231.109:3838/song_hotness_regression/

Choose the Variables for Regression for
Song Hotttness

☑ artistHotttness
☐ duration
☐ tempo
☐ songKey
☐ loudness
☐ timeSignature
☑ peakPosition
☐ songYear
☐ billboardYear
☐ mode

```
Call:
lm(formula = as.formula(paste("songHotttness ~", terms)), data = df)

Residuals:
     Min       1Q   Median       3Q      Max
-0.69133 -0.10586  0.02192  0.12318  0.46395

Coefficients:
                 Estimate Std. Error t value Pr(>|t|)
(Intercept)     4.136e-01  6.876e-03   60.16   <2e-16 ***
artistHotttness 5.344e-01  1.311e-02   40.77   <2e-16 ***
peakPosition   -9.584e-04  5.599e-05  -17.12   <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.1683 on 10271 degrees of freedom
  (1990344 observations deleted due to missingness)
Multiple R-squared:  0.1625,    Adjusted R-squared:  0.1624
F-statistic: 996.7 on 2 and 10271 DF,  p-value: < 2.2e-16
```
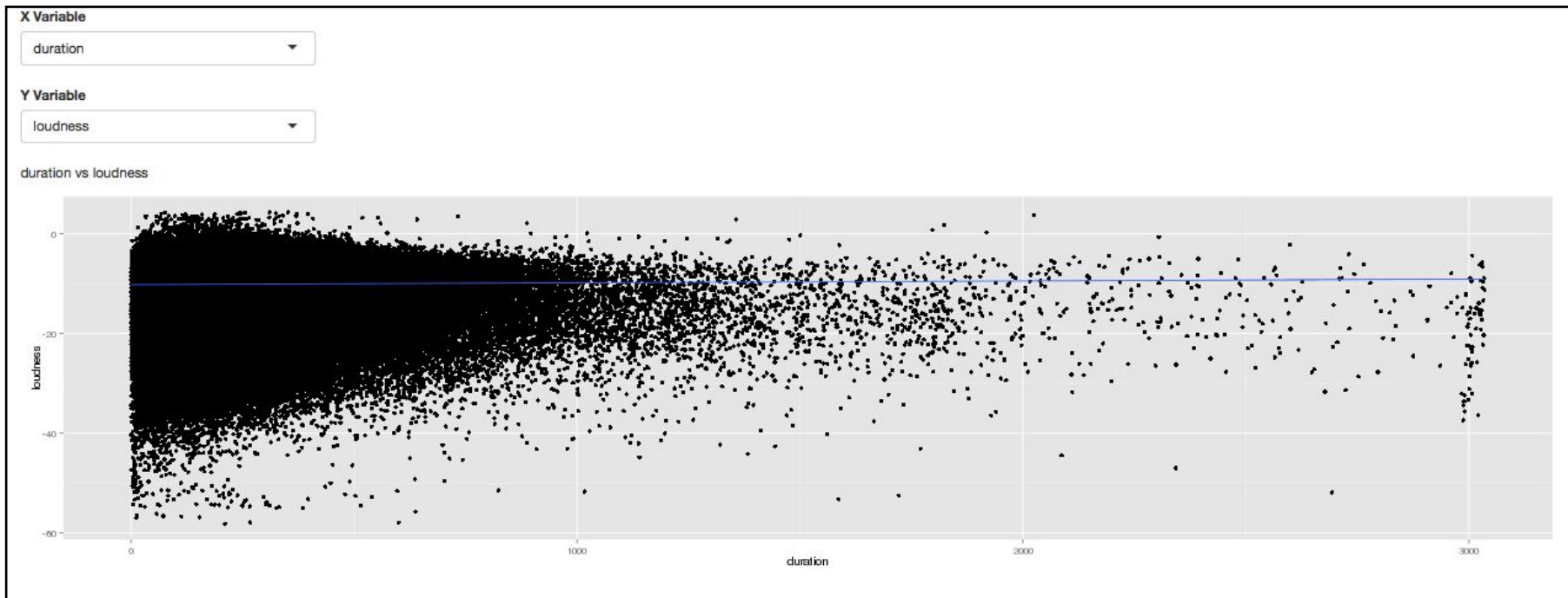
# Analysis: Shiny Example 3: Scatter Plots

# Architecture - Alternatives

- Processing Layer
  - Leave the single-song HDF5 files as-is
    - The "small files problem" wasn't nearly as problematic as the Hadoop/HDF5 problem.
  - Run the entire job in Spark
  - Convert HDF5 files to another format (e.g., JSON) from the outset
    - Human-readable
    - Hadoopable
    - Bigger data size, but ¯\\_(ツ)_/¯
- Serving Layer
  - DB - any relational database
  - BlinkDB - sampling database to speed up R visualizations
  - Tableau
  - Python with visualizations in D3

# Future of Project - Scaling

- As expected, the Hive portions of the job scale very well.
- HDF5 file processing was largely sequential.
  - Scaling would require devising a way to store HDF5 files in a data lake
  - Alternatively, they could be converted to a more compatible format at import time
- Shiny had difficulty with the volume of data
  - R loads all data points into RAM, which imposes scaling limits
  - Regression and plotting of such large volumes of data is slow
  - Could instead use a sampling database such as BlinkDB to control the amount of data