

CSCI 4253 Project Proposal

Version 3 - 12/4/18

Evan James

Mario Ramirez

Johnathan Kruse

Contents

Problem description	3
High-level solution architecture	4
Dataset	5
Challenges	7
Timeline	8
Checkpoint 1	9
Checkpoint 2	13
Sources	18

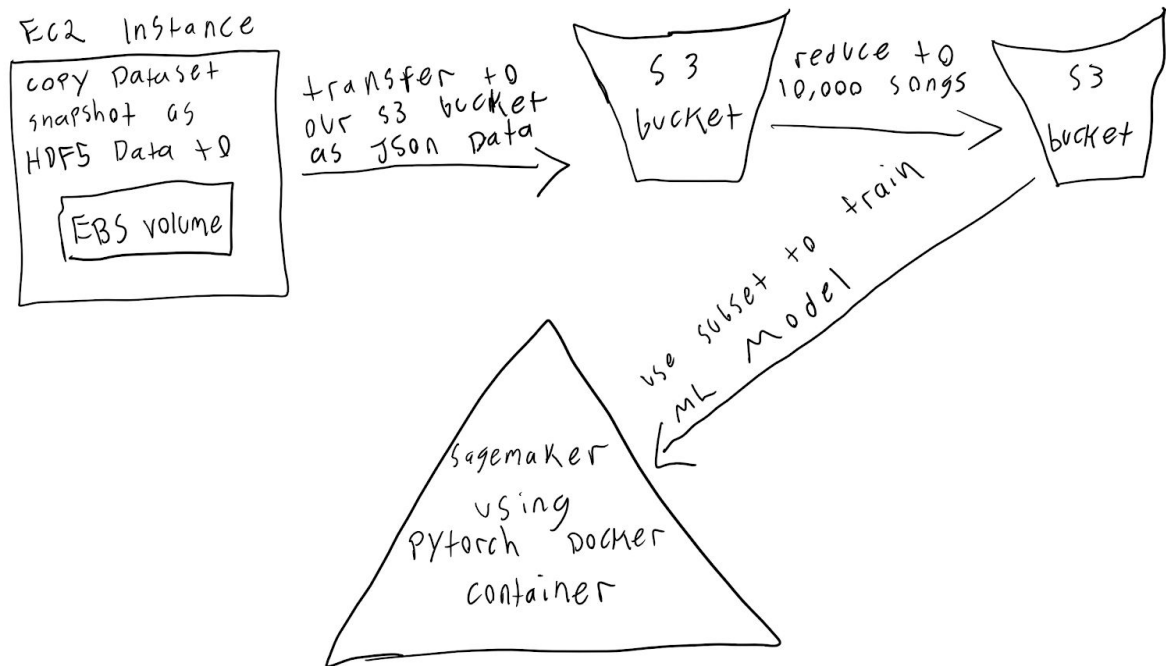
Problem description

Globally, the music industry is valued at approximately \$16 billion (Forbes, 2017). Furthermore, the number of songs downloaded and streamed each year is estimated to be well over a billion and is growing each year (Nielsen 2016 Music Year-End Report). As a musician, the market has become increasingly difficult to stay relevant. One of the greatest challenges is creating a song which will become “popular” and thus incur more consumption by listeners. Our goal is to help musicians around this problem by determining both if their song will be popular, and characteristics of what would make their song more popular.

For our project, we will utilize the Million Song Dataset -- found at <https://labrosa.ee.columbia.edu/millionsong/> -- to train a machine learning algorithm to predict whether songs will be popular or not. The MSDS contains over one million songs along with information of the song taken from The Echo Nest API. The song attributes include tempo, time signature, key, energy, danceability, and many others. The most interesting metric is “hottness” (note the spelling is intentional), which is a value which represents how popular the song was on a scale from 0 to 1. We will have a machine learning algorithm ingest the various attributes of each song and use it to be trained to guess whether any given song will be popular or not. Once the ML algorithm is trained, we will be able to query the Spotify API to get the necessary artifacts from a song and feed it into the algorithm to attempt to predict whether the song will be hot or not.

This is an interesting problem for many reasons. It can easily be seen how an algorithm like this would be useful in the music industry. It could be used to not only predict whether a song will be popular or not but could also be used to give insight into what kind of characteristics cause a song to be popular or not. The types of song attributes and their hottness can also be tracked over time, showing what kind of songs were popular each year and what kinds of trends happened over time with the sound of music. Having insight as to what makes music popular would be extremely useful for music streaming companies, such as Spotify, but also to record labels who have to decide on what kind of music to invest money and resources into.

High-level solution architecture



We plan to use Amazon SageMaker to implement the Machine Learning Model using the Pytorch framework. We will then want to use an s3 bucket that will contain the data from the million song dataset to train and test our Machine Learning Model on. SageMaker will be important for us to use because it can be directly connected to an s3 bucket which is where we will store our data. However, before we can implement this SageMaker ML Model to work with our s3 bucket we need to transform this data into a form that is more easily workable.

First, to get the raw data we will use the Amazon public dataset snapshot to copy the data into an EBS volume. This can then be copied into our s3 bucket which we will manipulate from the original HDFS format to JSON format. Once we have the data in a format we like we will reduce the size of the data to 10,000 songs to be more easily workable. After the data is set up correctly we can then connect our ML Model directly to our s3 bucket to start training and testing.

Dataset

For the project, we will use the Million Song Dataset (MSDS). The dataset contains 273 GB of song information compiled by Lab ROSA using data provided by The Echo Nest. The Echo Nest, a Spotify owned company, is a data collection organization whose focus is specifically on music. LabROSA is a Columbia University research organization whose focus is on intelligent sound recognition. LabROSA selected songs based on familiarity of the artist, song similarity and popularity. The dataset was last updated in November of 2015. Regarding overlap, the songs selected by LabROSA are unique (no duplicates). There are additional datasets provided by LabROSA that overlap and provide additional song data (e.g. lyrics). However, the current plan is to stick with the LabROSA Echo Nest data only.

LabROSA has stored the data in an HDF5 format. HDF5 is a file storage technology which allows very large heterogeneous data sets, file compression and search optimization. For each song, there are over fifty features (data about the song), all of which are not needed. Hence there will be preprocessing for eliminating irrelevant data features, which will also reduce the total file size. Depending on the overall reduced size, the processed data will be transferred to local storage (e.g. as JSON) to eliminate redundant preprocessing computation. Please see the table below for some sample fields.

Subset of Dataset Fields:

Field Name	Type	Description
Title	String	Name of the song
Loudness	Float	How loud the song is overall (dB).
Song Hottness	Float	An index rating between 0 and 1 (inclusive) of how popular a song was when released (e.g. 1 means very popular).
Tempo	Float	Song tempo (BPM).

More information can be found at: <https://labrosa.ee.columbia.edu/millionsong/pages/field-list>.

The MSDS is a static dataset, which can be accessed through a couple of different ways. The Open Science Data Cloud (OSDC) hosts the data as 26 tar files which can be downloaded locally. Alternatively, the data is available as an Amazon Public Dataset snapshot free of charge. All that is required is an Amazon Web Services (AWS) account which each member of the team already has access to. Therefore, we can retrieve the data by using an EC2 (Elastic Cloud Compute) virtual machine instance, attaching the data as AWS EBS (Elastic Block Storage) and then importing it.

During our initial development phase, we will store a small subset of the processed data (provided by LabROSA) as a CSV file (can import into HDFS too) in an S3 bucket. During our

final development phase, we will process the data, then store it in EBS or potentially an S3 bucket if costs exceed our expense limit. In both of these phases, the data will be imported in order for our machine learning algorithm to train on.

Challenges

The problem is difficult to solve because there is no straight-forward algorithm for making a “popular” song. There are many different factors that go into making a song and gauging the perception that will be generated by the public. For example, one factor that is hard to determine is that as time progresses the genre of “popular” songs changes with mainstream listeners. A song that may have been popular half a century ago would likely not be popular today.

Some of these factors may be too random to determine, in fact, while we might have a plethora of data at our disposal, we are not guaranteed to get any meaningful success from our machine learning algorithm. Regarding our solution, some of the main challenges will be evaluating, optimizing and selecting the best algorithm and accompanying parameters for our use case. Deciding which algorithm we select will impact our final results. Considering there are many algorithms at our disposal, this can be a challenging task. Some of the algorithms to consider are regression, neural networks, and more.

Furthermore, our group is not widely experienced with machine learning, especially at a large scale, so learning and understanding the concepts may be difficult. However, we will look to use libraries such as TensorFlow and other technologies to make this step easier. Additionally, each of these algorithms will have hyperparameters that will affect prediction accuracy. Learning the tuning process will be a challenge.

Timeline

Right now it is too early and we do not have enough information about Machine Learning Modeling to be able to effectively split up responsibilities. Machine Learning is something none of us have any experience in so that will be the main part that we all put our focus on. The other main part is how we will transform the HDF5 data from the million song dataset into JSON data to be used by our PYTorch ML Model. All of these parts are very new to us so we still need to dig deeper into these tools to find exactly how to implement them to properly split up the responsibilities.

For a timetable, we have a few goals that we want to meet before the due date and before the two checkpoints for the final project. Before checkpoint 1, we want to have the million song dataset reduced and transformed to JSON data so it is ready to be used by our ML Model. Also, we want our ML Model to be properly implemented in Amazon SageMaker so it is ready to be trained on our reduced dataset. Next, Before checkpoint 2, we want to have successfully trained our ML Model on the smaller, reduced dataset. Also, we want to have the larger subset of the data ready for use of training the ML Model. Finally, by the due date, we want to have successfully trained our ML Model on the entirety of the Million song dataset, or close to it, to effectively predict the potential “hottness” of a new song.

Checkpoint 1

Changes

No Changes to our initial outline.

Regarding the timeline, the original plan was to have the Echo Nest song dataset reduced into a CSV file full of approximately 10,000 samples with irrelevant columns removed. We also wanted to have our machine learning model ready to be trained for our next checkpoint. Both of these have been achieved, meaning we are still on track to meet our Checkpoint 2 goal of training our model on the sampled data points. Therefore our future timeline remains the same.

Also, we plan to add to the project by using the last.fm dataset determine which types of songs are hottest at certain times of the year. By doing this an artist can have data providing them with the appropriate time of the year to release their song that will increase their chance of it becoming a hot song based on some tags given to the song. The last.fm data set provides timestamps for songs when they were played along with some tags associated with the song.

Current Timeline

Work Done

For this checkpoint, two goals were completed:

1. Reduce the Echo Nest Dataset down to a smaller sample size. Keep necessary columns only.
2. Set up a prototype machine learning algorithm in AWS SageMaker, having selected the algorithm we wish to train on.

1. Reducing the dataset

The code to reduce the dataset is found in the root of the Github repo, titled `h5_to_csv.py`. The script is run with the syntax `python h5_to_csv.py <data directory>` where the data directory contains a series of sub-directories with the `.h5` files within. The Python script uses the `hdf5_getters.py` utility found in the `utils` folder, which makes the implementation significantly simpler. The `hdf5_getters.py` utility supplies functions which take an argument of the opened `.h5` file and returns the value in the field which the function was called to retrieve. For example, calling `hdf5_getters.get_artist_hottness(h5_file)` will return the artist hottness for the song in that `.h5` file.

Something to note is that the `h5_to_csv.py` file was tested on the 10,000 song subset of the MSDS, supplied from the MSDS website. All of the `.h5` files only have a single song in them, but it is possible that files in the entire dataset may contain more than one song per file. I don't think this will happen, but the possibility exists. In this case the `h5_to_csv.py` script could then be adapted to add functionality to iterate over the multiple songs in a file, since the `hdf5_getters.py` file contains a function to get the number of songs in a file.

The `h5_to_csv.py` script will iterate over every subdirectory and file in the input directory and output a file titled “`msds.csv`” which contains the information for every song found while traversing the directory. The data columns which are included are a significant subset of the entire possible columns, which brings the size of the data down considerably. This will be convenient moving forward, since we can take the entire ~250 GB million song dataset and make it a significantly smaller CSV file, which will be more efficient to work with. It also allows us to store less data, which will cost less to process on AWS.

2. Machine Learning

The work done for Machine Learning can be found in the Jupyter Notebook titled “`ml_prototype.ipynb`” (in the Github repo) which can be uploaded to AWS SageMaker. The notebook can quickly be configured to import a CSV file from local or AWS S3 into a pandas dataframe. Once the data is imported, the notebook can remove unnecessary columns and rows such as rows which are missing data. Note that data preprocessing is now completed in the previous step before being imported to the notebook, albeit when this notebook was first being developed it was necessary for testing the data as the previous preprocessing step was not complete at the time. It is still left in the notebook to make it easier for testing small samples of unprocessed data. Then the data is split into testing and training subsets.

Currently 80% of the data is used for training. Then the models are trained. Two models were tested, being K-Nearest-Neighbors (KNN) and Logistic Regression. The basic premise is that a song’s input features are:

- Duration
- Number of beats
- Loudness
- Tempo
- Mode
- Key
- Year

Then the model will predict if the song will be “hot” or “not hot” with a confidence between [0, 1]. We did consider using linear regression, albeit given that we wanted to use classification it felt more appropriate to Logistic Regression in its place. More details and reasoning can be found inside the Jupyter Notebook. After the models were trained, parameter tuning occurred. For example for KNN a list of different nearest neighbour values were used (e.g. 2, 5, 10) and then the testing accuracy was compared to find the best number of neighbors to use. Performance metrics such as prediction accuracy on test data were used to compare the models against each other. Notice it is important to understand that these findings (e.g. parameter tuning) can change once more data is imported for the next milestone. Currently our results indicate that KNN with approximately 6 or 7 neighbors would be the most appropriate.

To upload to SageMaker, simply go to AWS SageMaker through the user interface. Create a new notebook instance and then launch it. Once it is running, open the jupyter notebook, click

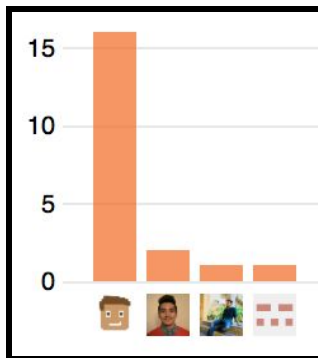
“Upload” in the top right corner and then upload the “ml_prototype.ipynb” notebook. An S3 bucket and CSV file must then be specified in order for the notebook to run properly.

Work To-Be Done

Now the model needs to train on a larger data set $\geq 10,000$ points to start out, at least for Checkpoint 2. Also by the next checkpoint we need to update our analysis that found KNN to be suitable as additional data could alter this conclusion. Lastly, by the end of the project all the data from Echo Nest needs to be pre-processed and trained on.

Last.fm data is already in json form and ready to be worked with so now we need to implement a way to grab the title of a song and map that with songs in the million song dataset where we will then look at the timestamps provided by the last.fm dataset.

Group Member Contributions:



Note: Some commits were larger and thus less frequent

Johnathan Kruse implemented the machine learning algorithms and comparison between models inside the entire Jupyter notebook. This is the content described underneath the 2. Machine Learning section just above.

Mario Ramirez researched ways to implement scikit-learn model on sagemaker, focusing on creating a docker container to run on sagemaker but ultimately Johnathan implemented using jupyter notebook. Researched ways to convert hdf5 dataset to csv or json. Evan implemented a way to transform to csv but I discovered a temporary solution in the form of a subset of million song dataset that was transformed to csv to be used for initial testing of ML model.

Evan James implemented a Python script to take the .h5 files of the MSDS and parse them into a CSV containing information for all the songs. The script only takes a subset of the available columns, since there is a lot of information which would not be useful when training the ML algorithm. The parser was originally written using the h5py library, and parsed the fields of each file. Further in development I found that there is a open source hdf5_getters.py utility provided by LabROSA.

Costs

Currently SageMaker has not exceeded \$5, albeit this price is likely to increase as larger datasets are imported. S3's cost has also been negligible to this point.

Dataset Management

One issue is some of the data (not a majority) is missing column values such as the year the song was made or the “hotness” score of the song. This data is removed though it should be noticed that this can decrease the accuracy of the model since now it has less data to train on.

Checkpoint 2

Meetings

Meeting 1:

- Time: 10:30am
- Date: 11/28/18
- Length: 20 min
- Participants: Evan, Mario, Johnathan (everyone)
- Topics covered:
 - LastFM data (has year of song) is in JSON
 - Will change it to CSV to make it easier to combine with million song dataset
 - Mario will work on this
 - Will explore data in Jupyter notebook to look for trends.
 - Evan will work on this
 - Will make a script to combine song years from LastFM csv file with Million song dataset
 - Johnathan will work on this
 - EDIT TO MEETING NOTES: Some of LastFM may already be incorporated to the Million Song Dataset
 - A more important priority is ensuring that we access more of the dataset and convert it to CSV per our original timeline. Currently we only have access to 10,000 songs.
 - Want to have ability to see top/hot songs for every year - later feature.
 - Eventually have ability to predict which year a song would be "hottest"
 - Have another check-in before turning in to see how everyone is doing.
 - Another meeting on the weekend (or weekday possibly) to start compiling

We then communicated via text messages, checking progress and making sure everything was completed on time.

Changes

No changes.

Per the last checkpoint it was decided that we would add an additional component to the project of combining the Million Song dataset with LastFM. This would give more information on the year of release for additional songs. We have converted the LastFM dataset from JSON to CSV to be able to work with it more easily. For our final turn-in we plan to have combined the datasets.

Current Timeline

Work Done

For this checkpoint, three tasks were completed:

1. Explore the dataset for more trends
2. Convert LastFM JSON to CSV
3. Get more data from the Million Song Dataset and upload to S3

1. Dataset Exploration -

Using a Jupyter Notebook and PyPlot, I implemented multiple visualizations of the data to get a better sense of the various trends. To begin with, rows with N/A values for the hottness are dropped. This is necessary to be able to properly classify the data and train the ML algorithm. Since we are training to be able to predict hottness, having rows of data without any value for hottness will be useless. Then the dataset is sorted to be in decreasing hottness.

Various visualizations are then created, showing the trend of tempo, loudness, and duration over time. Right now these graphs do not seem to give any good insight into the data, since it is only using the 10,000 song subset. Once the useless rows are dropped, that subset parses down to only about 4,000 datapoints. After the visualizations are implemented with the whole million song dataset, the visualizations might be more meaningful.

Aside from the visualizations, I implemented copying the snapshot of the MSDS into a EBS disk and mounting it onto an EC2 instance to run the parser on the full dataset. As of writing this, it is still in the middle of running the parser. Since the dataset is so large, I estimate the runtime will take around 16 hours.

2. LastFM JSON to CSV

The LastFM dataset is given in JSON format so in order to merge with the Million Song Dataset we needed to convert this data from JSON to CSV. First, we downloaded a subset of 10,000 songs from the LastFM dataset to begin conversion to CSV. In order to convert this subset a python script was developed to take in all 10,000 song data, convert to CSV, and converge to one single CSV file to be used with the Million Song Data Set.

The LastFM data comes with a lot of descriptors for each song including tags, timestamps, comparisons, etc. However, a lot of these descriptors we do not need so the python script also filtered each songs data to only take the song ID, song name, and timestamp. Only this data was needed to match songs to those in the Million Song Data Set and use the timestamp to determine best time to release a specific genre of song to maximize hottness.

3. More Data and Retraining

Note that there are specific instructions for the following under the parsing/ directory. Below is a higher-level description of what was completed.

Originally, for the first checkpoint a small subset of approximately 10,000 songs were used for the Machine Learning model to train on. This subset was a small sample of files in the HDF5 format provided by LabROSA on their website. Evan had made a python script for converting HDF5 to CSV for Checkpoint 1. For Checkpoint 2, our goal was to create a CSV file with more than the 10,000 songs given to us. We ended up creating a CSV file with almost 400,000 songs. One of the challenges with this was accessing this additional data. In order to obtain more HDF5 files, LabROSA provided the Million Song dataset on AWS as a public Elastic Block Storage (EBS) snapshot. An EBS snapshot is an exact copy of a chunk of storage at a certain “point-in-time” (per AWS documentation). **In order to use an EBS snapshot, it must be mounted as a volume for an EC2 instance.** A volume creates a replica of the snapshot and hosts it with our AWS account allowing an EC2 instance to mount it and make use of it.

For this process, we first looked up the EBS snapshot id found at <https://aws.amazon.com/datasets/million-song-dataset/>. We then went to EC2 and setup a new instance. Per normal configuration options, we set up a standard single md5.large instance with a security group that only involved only our local IP address. When prompted, we then attached a volume using the snapshot id from the url above. Note the volume must be at least equal in size to the snapshot, which was 500 GB. We then started the instance and used SSH to access it. From here we mounted the volume to a file path. We used the scp tool to transfer the Python script which can convert HDF5 to CSV. We ran for approximately 400,000 songs. This is about 40 times greater in size than the original sample subset. It is also important to recognize that the ELB snapshot is only available in US-East-1 (Virginia) which is where the EC2 instance must also be ran. The new larger file is stored on Google Drive and can be found under the data/msds_larger.md file.

Once the script finished running and gave a CSV file, we then uploaded it to an S3 bucket. From here we can now retrain our ML model per the original timeline. Currently, the ML model has been trained on James’ 10,000 song which had the optimal number of neighbours at ~6. For the final turn-in we will train on the new 400,000 song subset. We have also uploaded the data to Google Drive and there is a link to this under the data/ directory.

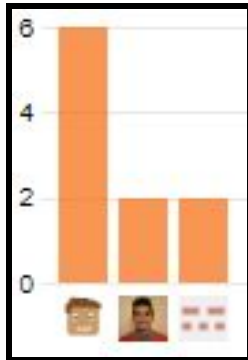
The steps described above were done mainly via the AWS web console and using SSH to access the EC2 instance. Inside the parsing/ directory there are more instructions on setting this up via web console and command line.

Work To-Be Done

Look into combing the Million Song dataset with the LastFM dataset to collect further information on the release year of songs. **We also have to retrain the ML model** on the new

larger set of data and also explore adding popularity by year such that our model can predict which year would have been the best to release a song.

Group Member Contributions:



Note: Some commits were larger and thus less frequent. Left to Right: Johnathan, Mario, Evan.

Johnathan Kruse got more data from the Million Song dataset and uploaded the file to S3 and Google drive. This was described under Part 3 of our goals completed inside the Timeline section above.

Mario Ramirez created the Python script that converted the LastFM JSON dataset to CSV format. This was described under Part 2 of our goals completed inside the Timeline section above.

Evan James implemented the dataset exploration Jupyter notebook looking for more trends in the data. This was described under Part 1 of our goals completed inside the Timeline section above.

Costs

Costs have increased due to the resources required for EC2 and volume mounting. We estimate that approximately \$20 have been spent. Going forward, costs will likely increase due to use of SageMaker and storing data in S3, **albeit we do not anticipate this to exceed this current cost** as we should not have to use EC2 and volume mounting going forward.

Dataset Management

Currently we now have access to a larger dataset stored as a CSV in an S3 bucket. The data is still cleaned inside the AWS SageMaker Jupyter notebook before being trained on. Given the size of the dataset, we would like to clean the data during the “combining phase” of the Million Song dataset with LastFM in order to reduce the required storage size.

Challenges

Challenges have mainly been with getting more data from the Million Song Dataset as the process was not initially easy to determine. However, we have now determined how to get more data, albeit it can be a slow process.

i'm a bit worried that y'all are going to be in trouble for the final report. the progress here didn't look great, but i might be missing something. please continue working hard, and don't hesitate to come meet with me if you are having any issues. good luck!

10/10

Sources

Million Song DataSet:

Thierry Bertin-Mahieux, Daniel P.W. Ellis, Brian Whitman, and Paul Lamere.
The Million Song Dataset. In Proceedings of the 12th International Society
for Music Information Retrieval Conference (ISMIR 2011), 2011.

LabRosa:

Million Song Dataset, official website by Thierry Bertin-Mahieux,
available at: <http://labrosa.ee.columbia.edu/millionsong/>

The Echo Nest <http://the.echonest.com/>

HDF5 <https://support.hdfgroup.org/HDF5/whatishdf5.html>

AWS <https://aws.amazon.com/datasets/million-song-dataset/>

Forbes Value of the Music Industry:

<https://www.forbes.com/sites/hughmcintyre/2017/04/25/the-global-music-industry-grew-by-6-in-2016-signalling-brighter-days-ahead/#5dc3a83163e3>

Nielsen 2016 Music Year-End Report:

<https://www.nielsen.com/us/en/insights/reports/2017/2016-music-us-year-end-report.html>