# Rap Recommender

**Building A Hybrid Recommendation Engine From Scratch**

Owen Curtis  DSI-9

# Goals

**Build Rap Recommender** + **Execute Learning Agenda** + **Publish Dataset**

# The Data

**800** Rappers/Artists

**25,000** Tracks
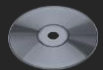
**17,000** Album Reviews

**4MM+** Words

**211** Hours of Audio

**151,000** Users

# The Data

**800** Rappers/Artists

**MM+** Words

**25,0** ...

f Audio

**17,** ...

**151,000** Users

**80%** of tracks
fucking profane!

# The Process

Source ➤ Collect ➤ Engineer ➤ Explore ➤ Recommend ➤ Combine

# The Learning Agenda 🧪

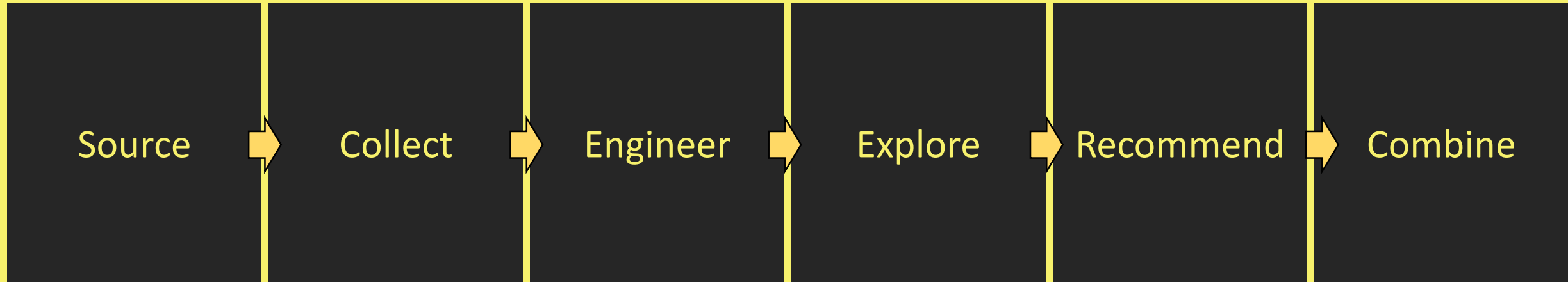| Source | ➤ | Collect | ➤ | Engineer | ➤ | Explore | ➤ | Recommend | ➤ | Combine |
|--------|---|---------|---|----------|---|---------|---|-----------|---|---------|

**Scraping:** Is Selenium an effective library for more bespoke web scraping needs?

**Matching:** What are best practices for fuzzy matching on similar strings?

**Insights**: How has the genre changed? Is hip hop dead as an art form?

**Topics:** What is the most effective way to extract topics from text data?

**Production:** Are audio features an effective means of predicting song producers?

**Recommender:**
- How does the output differ between approaches (collab v content)?

- What are the pitfalls of recommender systems? Which do we encounter?

- What is the appropriate approach for a Hybrid recommender?

**Source** ➤ Collect ▸ Engineer ▸ Explore ▸ Recommend ▸ Combine

- **Audio**
- **Reviews**
- **Lyrics**
- **User Data**

# Audio Data

# Review Data

# Review Data

# Lyrics & Users

Source → **Collect** → Engineer → Explore → Recommend → Combine

- **Querying/Scraping**
- **Challenges**
- **Cleaning**
- **Fuzzy Matching**

# **Collect:** An Iterative Process

Refine

Initial
Artist List

Refine

Manual Checks
Re-pulls
String Cleaning
Search Based

Refine

Refine

GENIUS

Refine

YouTube

# **Collect:** Approach

| Audio | Reviews | Lyrics & Users |
|---|---|---|

**Audio**



- Pre-existing API python wrappers

**Reviews**



- BeautifulSoup



- Kaggle / Google Datasets
- Pre-existing API python wrappers

**Lyrics & Users**



- **Lyrics:** Pre-existing API python wrappers

- **Users:** a bit more complicated...
  - No API Endpoint
  - Infinite Scroll
  - Large scale scrape

# **Collect:** User Data

# **Collect:** User Data

| Pros |
|------|
| • High degree of customization |
| • If it's in the HTML, you can get it |
| • Safeguards in place for interruptions |

| Cons |
|------|
| • Scalability/Speed |
| • Ads / Scripts |

**VS**

| Pros |
|------|
| • Well structured data |
| • Speed/Scalability |
| • No issues with ads/scripts |

| Cons |
|------|
| • Limited to data provided |

**Learning**
- Selenium is an excellent tool for things like automated site testing
- An effective tool for scraping when there are no other options (ex form filling, "clicking")

# We have our data. Let's check it out!!!!

| Source | Artists | Album | Track |
|---|---|---|---|
| (Spotify logo) | A$AP Rocky ft. Black Hippie | Long Live A$AP | U.E.N.O |

# Oh....Oh god....

| Source | Artists | Album | Track |
|---|---|---|---|
| (Spotify) | A$AP Rocky ft. Black Hippie | Long Live A$AP | U.E.N.O |
| (Spotify) | A$AP Rocky | Long Live A$AP | U.E.N.O. (ft. Jay Rock, Ab Soul, Schoolboy Q) |
| (GENIUS) | A$AP Rocky ft Jay Rock, Ab Soul, Schoolboy Q | Long.live.A$AP | U.E.N.O. |
| (R) | A$AP Rocky | Â°Ã¬â‚¬LongÂ°Ã¬â‚¬liveÂ°Ã¬â‚¬ AP | |

Time to switch projects?

# Collect: Preprocessing

| Artist Match Rate | Album Match Rate | Track Match Rate |
|:---:|:---:|:---:|
| **73%** | **54%** | **38%** |

"EP"
"LP"
"Edition"
**Text Reduction**
*("Volume" → "Vol" → "V")*

**Stripping**
Parentheticals, Bracketed Info

Punctuation     Special Chars     De-concatenation     Spacing     Gerunds

| **73%** | **62%** | **45%** |
|:---:|:---:|:---:|

# Collect: Cleaning

## Missing Data

| 7,000 Dates | 1,000 Albums | 6,000 Producers |
|:---:|:---:|:---:|
| *Imputed with mean* | *Imputed where possible* | *TBD (See Next Section)* |

# **Collect:** Fuzzy Matching?

**Matching:** What are best practices for fuzzy matching on similar strings?

| Jaro-Winkler Distance |
|---|
| Hamming Distance |
| Dice Distance |
| **Levenshtein Distance** |
| q - gram Distance |

*Levenshtein distance (or edit distance) between two strings is the number of deletions, insertions, or substitutions required to transform source string into target string.*

# A$AP Rockey | U.E.N.O

# ASAP Rocky | U.E.N.O.

# Distance: 3

# Collect: Fuzzy Matching?

**GENIUS**

**Matching:** What are best practices for fuzzy matching on similar strings?

| Model Features |
| --- |
| Jaro-Winkler |
| Hamming Distance |
| Dice Distnace |
| Levenshtein Distance |
| q - gram Distance |
| String Length |

**GENIUS**

Supervised Classification

N=1,500

## Performance

■ Train Acc   ■ Test Acc

*Baseline: 53%*

Logistic Regression   MultinomialNB   RandomForest   AdaBoost   GradientBoost

# Collect: Fuzzy Matching?

**GENIUS**

| | genius_artist_album | spotify_artist_album | lev | match | str_len_genius | str_len_difference | ham | dice | jaro_wink | probs |
|---|---|---|---|---|---|---|---|---|---|---|
| 380 | dj drama gangsta grillz: the album | dj drama gangsta grillz the album | 99 | 1 | 34 | 1 | 10 | 0.95082 | 0.99 | 0.980650 |
| 101 | spose the audacity | spose the audacity ( ) | 100 | 1 | 18 | 4 | 0 | 0.888889 | 0.96 | 0.976050 |
| 102 | spose happy medium | spose happy medium ( ) | 100 | 1 | 18 | 4 | 0 | 0.894737 | 0.96 | 0.976050 |
| 4 | open mike eagle rappers will die of natural ca... | open mike eagle rappers will die of natural ca... | 100 | 1 | 50 | 3 | 0 | 0.977778 | 0.99 | 0.972768 |
| 104 | spose preposterously dank | spose preposterously dank () | 100 | 1 | 25 | 3 | 0 | 0.93617 | 0.98 | 0.970815 |

| Learning | • There are several distance metrics which can be leveraged for text processing; selection context dependent<br>• These metrics can be fed in as features for supervised classification; high sample recommended |
|---|---|

# **Collect:** Fuzzy Matching?

**Matching:** What are best practices for fuzzy matching on similar strings?



**Probable Matches**     +     **Crowdsourced Validation**

**Learning**
- There are several distance metrics which can be leveraged for text processing; selection context dependent
- These metrics can be fed in as features for supervised classification; high sample recommended

Source → Collect → **Engineer** → Explore → Recommend → Combine

- **Reviews**
- **Audio**
- **Lyrics**

RAPREVIEWS    YouTube

"hey guys Anthony Fantana here Internet's busiest music nerd this review is a little late but you know it's because I spent all night last night compiling my favorite music videos EPS album's tracks singles I think I'm going to have a pretty good set of videos coming up toward the end of this year just a best-of thing alright we have a long long review this time around and it's because you know what we have to rap about rap there's a growing sentiment in some of my videos that I don't like rap that I don't like hip-hop and of course this foul stench is emanating mostly from my kanye and my big boy reviews people get hung up on the scores don't actually listen to what I say what can you expect but if you did hear what I said in those reviews you would know that some of my favorite songs of this year have come from those albums nevermind the fact that most people gloss over my roots review and tie early hip hop has really been kind of weird for me this year a lot of the major label releases the big ones have been somewhat decent for the most part but it's not like there haven't been hyped artists that I didn't really care for like all the weed wrap on currencies new LPS kind of makes it boring for me and big crits debut was alright and I have no clue why some of my videos are getting waka flocka flame requests there have been indie label releases that I liked for example gangrene or strong arms steady but for the most part Stones Throw rhyme Sayers Def Jux these indie hip hop labels have not been bringing the heat they haven't been bringing the competition to the majors…."

# Engineer: Reviews



Pre-Processing

Combined Review Text

Tfidf Vectorization

# **Engineer:** Audio

**Spotify**

**Custom**

# Engineer: Audio

## Spotify

| | | |
|---|---|---|
| Acousticness | Danceability | Energy |
| Instrumentalness | Liveness | Loudness |
| Speechiness | Valence | Tempo |

# Engineer: Audio

API → 30 second preview URLs → **Librosa** (aws)

| | |
|---|---|
| Tempo | Rolloff |
| Chroma | ZCR |
| Spectral Center | MCR |
| Spectral BW | Spectograms* |

*https://towardsdatascience.com/music-genre-classification-with-python-c714d032f0d8

# Engineer: Audio

**Let's look at a few examples…**

|  | Tempo | Chroma | Zero Crossing Rate |
|---|---|---|---|
| **Description** | Beats per Minute | Representation of Musical signals in terms of octaves. | Rate at which music signal changes from positive to negative, or vice versa |
| **Use Case** | High Pace, Low Pace | Identifying the tones of a song | Percussion |

*https://towardsdatascience.com/music-genre-classification-with-python-c714d032f0d8

# Engineer: Audio

We're still missing producer data on ~25% of tracks.

Can we use our audio features to identify them?

# Engineer: Audio

**Question:** Are audio features an effective means of predicting the track producer?

## Neural Network Setup

| Input Layer (16) |
| Hidden Layer (32, L2 Reg = 0.1) |
| Hidden Layer (32, L2 Reg = 0.1) |
| Output ('Softmax') |

*Epochs: 1000, Batch size: 128*
*N = 3,000*

## Accuracy

| Train | Test |
|-------|------|
| ~33% | ~28% |

| Baseline |
|----------|
| 4% |

## Top 5 Accuracy



*https://towardsdatascience.com/music-genre-classification-with-python-c714d032f0d8

# Engineer: Audio



**Question:** Are audio features an effective means of predicting the track producer?

**Speechiness**

**Unique Word Count**

# Engineer: Audio

**Question:** Are audio features an effective means of predicting the track producer?



EXAMPLE MIX PLACEMENT

50% LEFT — 50% RIGHT

HI HATS
TOM 2
BASS
LEAD GUITAR   SNARE
KICK
PIANO   ACOUSTIC
TOM 1   YOU   TOM 3
LEAD VOCAL
BG VOCAL 1
RHYTHM GUITAR   KEYS - PADS
VOLUME   VOLUME
NECESSARY?
BG VOCAL 2
BG VOCAL 3

In stereo, vocals recorded at "center" of music mix, can be isolated

Original     Reduced Vocals

# Engineer: Audio

**Better Processing**

**+**

**Audio Imaging**

**+**

**Convolutional Neural Network**



Mel-frequency spectrogram

**Learning**

- Our current set of features are not sufficient to predict a track's producer in order to impute
- Consider Convolution Neural Network

GENIUS

# **Engineer:** Lyrics

**GENIUS**

---

## **Rhyme Density**

### *A/B/A/B Scheme*

Here it is the groove slightly transformed
Just a bit of a break from the norm
Just a little somethin' to break the monotony
Of all that hardcore dance that has gotten to be
A little bit out of control it's cool to dance
And think of the summers of the past
Adjust the base and let the alpine blast
Pop in my CD and let me run a rhyme
And put your car on cruise and lay back 'cause this is summertime

*Will Smith - Summertime*

### *Internal Rhyming*

I'll knock you bitches into next week with a haymaker,
And straight-razor your face when you land seven days later,
I said pray your soul to keep when you go to sleep,
But you sold the lease on your own beliefs like Roman priests,
You probably told the beast Hip-Hop needs it's own police,
To patrol the streets and shows whenever something dope's released…

*Diabolic – Stand By*

# **Engineer:** Lyrics

## **Rhyme Density**

### **Phonetic Representation**

| but | you | sold | the | lease |
|---|---|---|---|---|
| b ʌ | j uː | s **əʊ** l d | ð ə | l **iː** s |

| on | your | own | beliefs | like |
|---|---|---|---|---|
| ɒ n | ɔː ə | **əʊ** n | b ɪ l **iː** f s | l aɪ k |

| Roman | | priests |
|---|---|---|
| r **əʊ** m ə n | | p r **iː** s t s |

$$\text{Avg Rhyme Density} = \frac{\text{\# of Potential Rhymes Per Line}}{\text{\# of Lines in Song}}$$

...
I'll knock you bitches into next week with a haymaker,
And straight-razor your face when you land seven days later,
**I said pray your soul to keep when you go to sleep,**
But you sold the lease on your own beliefs like Roman priests,
You probably told the beast Hip-Hop needs it's own police,
To patrol the streets and shows whenever something dope's released,
...

# Engineer: Lyrics

**GENIUS**

## Vocabulary

### Complexity
(# of 3+ syllable words in vocab)

### Size
(# of unique words in first 5K words)

## Sentiment (VADER)

### Positive

### Negative

### Neutral

**NOTE:** VADER is Lexicon based. Can get better results using supervised approach

**LDA**
Latent
Dirichlet
Allocation

# **Engineer:** Lyrics

GENIUS

- Remove Stopwords
- **Profanity**
- Lemmatize
- Bag of Words

- Build Your Dict
- Select # Topics
- Score Model

# **Engineer:** Lyrics

**Question:** What is the most effective way to extract **topics** from text?

- **Coherence Score:** Degree of semantic similarity between top words used to describe a topic

- Must be balanced with interpretability...



**Range: 0** (Bad) to **1** (Perfect)

# Engineer: Lyrics

GENIUS

```
[(0,
  '0.114*"get" + 0.065*"money" + 0.020*"gon" + 0.019*"ride" + 0.016*"talk" + '
  '0.016*"tryna" + 0.016*"bout" + 0.014*"work" + 0.013*"diamond" + '
  '0.013*"buy"'),
 (1,
  '0.071*"get" + 0.021*"mother" + 0.021*"keep" + 0.020*"see" + 0.019*"boy" + '
  '0.019*"come" + 0.017*"real" + 0.016*"run" + 0.016*"go" + 0.015*"hit"'),
 (2,
  '0.099*"go" + 0.072*"know" + 0.058*"let" + 0.044*"baby" + 0.041*"want" + '
  '0.038*"get" + 0.034*"come" + 0.033*"girl" + 0.029*"make" + 0.028*"back"'),
 (3,
  '0.134*"bad" + 0.110*"stop" + 0.066*"drop" + 0.059*"ready" + 0.058*"pop" + '
  '0.057*"top" + 0.036*"notil" + 0.023*"party" + 0.013*"business" + '
  '0.013*"click"'),
 (4,
  '0.048*"life" + 0.036*"time" + 0.029*"take" + 0.026*"live" + 0.022*"see" + '
  '0.022*"go" + 0.019*"mind" + 0.018*"lose" + 0.017*"find" + 0.015*"world"'),
 (5,
  '0.045*"say" + 0.039*"know" + 0.034*"get" + 0.026*"tell" + 0.020*"think" + '
  '0.020*"really" + 0.020*"never" + 0.018*"man" + 0.017*"even" + 0.016*"see"'),
 (6,
  '0.055*"get" + 0.032*"put" + 0.023*"high" + 0.023*"look" + 0.018*"smoke" + '
  '0.018*"hand" + 0.017*"roll" + 0.014*"body" + 0.013*"drink" + 0.013*"light"'),
 (7,
  '0.021*"make" + 0.013*"p" + 0.010*"back" + 0.008*"rap" + 0.008*"come" + '
  '0.007*"flow" + 0.007*"sound" + 0.007*"blood" + 0.007*"beat" + '
  '0.007*"rhyme"'),
 (8,
  '0.457*"be" + 0.098*"s" + 0.057*"can" + 0.030*"ill" + 0.026*"will" + '
  '0.024*"would" + 0.022*"have" + 0.010*"going to" + 0.009*"smokin" + '
  '0.007*"imma"'),
 (9,
  '0.169*"love" + 0.086*"feel" + 0.060*"never" + 0.045*"leave" + 0.034*"away" '
  '+ 0.032*"fall" + 0.027*"heart" + 0.025*"make" + 0.023*"break" + '
  '0.021*"hurt"')]
```

| Topics |
| --- |
| *The Hustle* |
| *Family* |
| *Lust* |
| *Partying* |
| *Reflection* |
| *Storytelling* |
| *Drugs* |
| *Skills* |
| *Aspirations* |
| *Love* |

**Learning**
- LDA is a powerful tool for topic extraction but will require subjectivity (num topics, topic interpretation)
- Can be used to programmatically generate tags based on content

Source → Collect → Engineer → **Explore** → Recommend → Combine

- **Lyrical Deep Dive**

# Lyrical Insights: Vocab

GENIUS

## Rappers by Vocabulary

### Artist Vocab Size and Complexity



### High Vocab Complexity

| | artist_clean | artist_vocab_complexity |
|---|---|---|
| 13563 | loonie | 0.420962 |
| 1650 | bewhy | 0.316637 |
| 23563 | wc no beat | 0.316547 |
| 936 | amill leonardo | 0.315754 |
| 3737 | ceza | 0.301329 |
| 11015 | k.a.a.n. | 0.301186 |
| 1618 | becky g | 0.253669 |
| 20355 | solillaquists of sound | 0.244256 |
| 1521 | bahamadia | 0.243795 |
| 6135 | doseone | 0.243747 |

### Low Vocab Complexity

| | artist_clean | artist_vocab_complexity |
|---|---|---|
| 7171 | fatman scoop | 0.044776 |
| 7100 | fat pat | 0.044248 |
| 15355 | mitchy slick | 0.042373 |
| 1340 | asian da brat | 0.038567 |
| 15113 | mike jones | 0.037975 |
| 7205 | fetty wap | 0.035363 |
| 16335 | nle choppa | 0.034483 |
| 19423 | scrilla | 0.032520 |
| 12921 | lil mosey | 0.028571 |
| 14772 | mc mong | 0.000000 |

# Lyrical Insights: Flow

# Lyrical Insights: Complexity

GENIUS

| Most Complex Track | Least Complex Track |
|---|---|

**Immortal Technique:** Speak Your Mind

**JPEGMAFIA:** *My Thoughts On Neogaf Dying*



*Based on track vocab complexity



*Based on track vocab complexity

# Lyrical Insights: Complexity

| Most Complex Track | Least Complex Track |
|---|---|
| **Immortal Technique:** Speak Your Mind | **JPEGMAFIA:** *My Thoughts On Neogaf Dying* |



Only a fucking imbecile would think their uncorrectable
Cause your susceptible to becoming more than a spectacle
Remember that your flesh, and blood and your bodies dissectible
I'll beat you until your a vegetable
And wake up in a hospital covered in poisonous chemicals
In a fetal position wit your face sewn to your testicles
Thinking that you were kidnapped by extraterrestrials
You got heart? I'm the blood that pumps in your ventricles
Technique, I'm like your soul nigga.. indispensable

# Lyrical Insights: Sentiment

GENIUS

☺

High Positive Sentiment Score

☹

High Negative Sentiment Score

**CORRECT**

**Your Love**
**Mick Jenkins**

Produced by KAYTRANADA
Album Wave[s]

**Hate**
**Vado**

Featuring Coko
Produced by Lee On The Beats
Album Sinatra 2.5

**INCORRECT**

**Wow**
**21 Savage**

Produced by Sonny Digital
Album Slaughter King

**Wicked**
**Future**

Produced by Metro Boomin & Southside
Album Purple Reign

# Lyrical Insights: Macro Changes

GENIUS

**Question:** How has the genre changed? Is hip hop dead as an art form?



Hip Hop is Dead and This Generation Killed it

Brian Brewington [Follow]
Nov 25, 2017 · 4 min read ★



## Hip Hop

## 1979 - 2017

# Lyrical Insights: Macro Changes

GENIUS

**Question:** How has the genre changed? Is hip hop dead as an art form?

Hip Hop is Dead and This Generation Killed it

Brian Brewington  Follow
Nov 25, 20

HIPHOP DX

f 1.6M    614K    323K    634K

HOME    NEWS    SINGLES    VIDEOS    REVIEWS    EDITORIALS    RELEASE

**5 Things That Killed Hip Hop**

February 20, 2007 | 12:00 AM
by J-23

25% OFF    ALL INSOLE
FREE SHIPP
Offer expires 12/

SHOP NOW    SUPERfe

HIP HOP

RECENT EDITORIALS

Hip Hop

1979 - ~~2017~~

2007*

# Lyrical Insights: Macro Changes

GENIUS

Hip Hop is Dead and This Generation Killed it

**Is Hip-Hop Dead?**

Follow

RomeluTrill-Kaku 5,915

▲ -4 ▼

When the term was first coined by Nas it was during a time(2006) where true lyricism in mainstream hip hop was lacking. There was a movement towards more party oriented songs and crunk music was in it's heyday. However, how relevant is that

HIP HOP

RECENT EDITORIALS

SHOP NOW   SUPERfe

Hip Hop

1979 - ~~2017~~

~~2007*~~

2006*

# Lyrical Insights: Macro Changes



**Question:** How has the genre changed? Is hip hop dead as an art form?

**Topics**

**Flow and Vocab**

Rhyme Density and Vocab Complexity Over Time

**2015: Migos** first major mixtape

# Lyrical Insights: Macro Changes

GENIUS

> **Question:** How has the genre changed? Is hip hop dead as an art form?

**Google Trends: 'Mumble Rap' Search Volume**

Oct 2018

July 2016

| | |
|---|---|
| 100 | |
| 75 | |
| 50 | |
| 25 | |
| Jan 1, 2004 | Oct 1, 2008    Note    Jul 1, 2013    Apr 1, 2018 |

**Learning**
- While it's not dead, it's become much more about flow than about the actual lyrical content

# Recommend: Data Summary

**Track**
38 Features

| Track | Album | Artist |
|---|---|---|

**Lyrical**

| Rhyme Density | Vocab Size | ... | Sentiment | Topics |
|---|---|---|---|---|

**Audio**

| Tempo | Speechiness | ... | ZCR | Rolloff |
|---|---|---|---|---|

**Other**

| Date | # Follows | Popularity |
|---|---|---|

# Recommend: Data Summary

**Track**
38 Features

|  | | | | Lyrical | | | | | Audio | | | | | Other | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| Track | Album | Artist | | Rhyme Density | Vocab Size | … | Sentiment | Topics | Tempo | Speechiness | … | ZCR | Rolloff | Date | # Follows | Popularity |

**Reviews**
1000 Words

|  | | | | Tfidf Word Freqs | | | | | | | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| Track | Album | Artist | | Word 1 | Word 2 | Word 3 | Word 4 | Word 5 | Word 6 | Word 7 | Word 8 | Word 9 | Word 10 | … | Word 1000 |

**Followers**
151K

|  | | | | Binary Follow Flags | | | | | | | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| Track | Album | Artist | | N00b_saibot | XXX COD420 | Gooch Crucible | Squiggly Bob | $mokeyyy | kartashov | TruSwag | YBN Mystique | Astroboy | efray | … | Squiggly Bob |

# Recommend: Types of Recommenders

| Bucket | Requirements | Advantages | Disadvantages |
|---|---|---|---|
| Knowledge-Based | • Significant knowledge of the domain in question | • Minimal reliance on data or specific user preferences.<br>• Approach for high complexity/barrier items | • Need extensive domain knowledge<br>• Scalability a factor |
| Heirarchical | • Clean product/content management system | • Simple to implement (I bought a printer, so I probably need ink) | • Limited personalization opportunity |
| ★ Content Based Filtering | • Clean and consistent metadata at the item level | • No community required | • Overengineering possible, leading to **unexpected** results<br>• Best when aligned to user preferences |
| ★ Collaborative Filtering | • Large historical data set<br>• Set of users and preferences | • No domain knowledge required<br>• Simple implementations can be effective | • Sparsity of Data<br>• Cold-Start User Problem |

# Recommend: Watchouts

| Watchouts | Description | Example |
|---|---|---|
| **Synonymy** | • Products or content that looks **very similar** but is in fact different. | • Eminem – *Without Me*<br>• Emimen – *Without Me* (Remastered)<br>• Eminem – *Without Me* (Radio Edit) |
| **Novelty vs Correctness** | • Combining novelty and correctness produces serendipitous moments for the user | • Danny Brown – *Bruiser Brigade*<br>• Casey Veggies – *PNCINTLOFWGKTA* |
| **Scalability** | • How effectively an engine can be scaled as underlying data grows<br>• *Sparsity* can help here | • Collaborative filter + massive social networks<br>• Page Views |
| **New Additions** | • The introduction of a new user or product into your algorithm – how will they be treated? | • New user with no profile<br>• Users with low rating volume<br>• New item with minimal metadata |

# Recommend: Measurement (Cosine Similarity)

# Recommend: Measurement (Cosine Similarity)



Close Relationship

RECOMMEND

No Relationship

~~RECOMMEND~~

# Recommend: Measurement (Macro KPIs)

# **Recommend:** Collaborative Filtering (Followers)

| Current Artist | Users Also Listen To… | Similarity |
|---|---|---|
| Kendrick Lamar | Schoolboy Q | 0.38 |
| | A$AP Rocky | 0.37 |
| | Nas | 0.34 |
| | J. Cole | 0.33 |
| | Chance the Rapper | 0.31 |
| | Big Sean | 0.30 |
| | Ab-Soul | 0.30 |
| | … | |

| Current Artist | Users Also Listen To… | Similarity |
|---|---|---|
| Chance The Rapper | Childish Gambino | 0.36 |
| | Kendrick Lamar | 0.33 |
| | Schoolboy Q | 0.31 |
| | Ab-Soul | 0.29 |
| | Mac Miller | 0.28 |
| | A$SAP Rocky | 0.27 |
| | Vic Mensa | 0.24 |
| | … | |

**Diversity**    Serendipity    Novelty    **Relevancy**

# Recommend: Content Based (Review Text)

| Current Album | Related Reviews | Similarity |
|---|---|---|
| | Kendrick Lamar: *DAMN* | 0.89 |
| | Kendrick Lamar: *Good Kid, MAAD City* | 0.78 |
| | Kendrick Lamar: *Section .80* | 0.75 |
| | Kendrick Lamar: *Untitled* | 0.73 |
| | Ab-Soul: *These Days* | 0.66 |
| Kendrick Lamar *To Pimp A Butterfly* | Logic: *The Incredible True Story* | 0.54 |
| | Schoolboy Q: *Habits and Contradictions* | 0.53 |
| | ... | |

| Current Album | Related Reviews | Similarity |
|---|---|---|
| | Chance the Rapper: *The Big Day* | 0.86 |
| | Chance the Rapper: *Coloring Book* | 0.75 |
| | **Donnie Trumpet : *Surf*** | 0.64 |
| | ~~Jeru The Demaja: *Divine Design*~~ | 0.61 |
| | Childish Gambino: *Because the Internet* | 0.57 |
| Chance The Rapper *Acid Rap* | Mackelmore: *This Unruly Mess I've Made* | 0.55 |
| | **Bobby Womack: *The Bravest Man In...*** | 0.54 |
| | ... | |

Diversity   Serendipity   Novelty   Relevancy

# **Recommend:** Content Based Filtering (Track Features)

| Current Track | Current Album | Current Artist |
|---|---|---|
| *Alright* | To Pimp a Butterfly | Kendrick Lamar |

| Track | Album | Artist | Similarity |
|---|---|---|---|
| Lost Ones | *Cole World* | J Cole | 0.32 |
| ~~Kiss Kiss~~ | ~~*Exclusive*~~ | ~~Chris Brown~~ | ~~0.31~~ |
| Too Good | *Views From the 6* | Drake | 0.31 |
| Can't Feel My Face | *Beauty Behind Madness* | The Weeknd | 0.29 |
| **Land of the Snakes** | ***Born Sinner*** | **J Cole** | **0.28** |
| sing about me im dying | *To Pimp a Butterfly* | Kendrick Lamar | 0.27 |
| Cocoa Butter Kisses | *Acid Rap* | Chance the Rapper | 0.27 |
| **What in XXXTarnation** | ***Members Only*** | **XXXTentacion** | **0.26** |

Diversity   Serendipity   Novelty   Relevancy

# Recommend: Content Based Filtering (Track Features)

| Current Track | Current Album | Current Artist |
|---|---|---|
| *Favorite Song* | Acid Rap | Chance the Rapper |

| Track | Album | Artist | Similarity |
|---|---|---|---|
| Global | *Street Gossip* | Lil Baby | 0.41 |
| ~~Codeine Dreaming~~ | ~~*Project Baby 2*~~ | ~~Kodack Black~~ | ~~0.40~~ |
| Hot Shower | *The Big Day* | Chance The Rapper | 0.40 |
| I'm Straight | *Harder Than Ever* | Lil Baby | 0.38 |
| This is a Hit | *Beautiful Loser* | Kyle | 0.33 |
| ~~The Rain~~ | ~~*Supa Dupa Fly*~~ | ~~Missy Elliot~~ | ~~0.29~~ |
| ~~Gunwalk~~ | ~~*I Am Not Human Being2*~~ | ~~Lil Wayne~~ | ~~0.29~~ |
| **Giv no Fucks** | ***Late Nights*** | **Jeremih** | **0.28** |

Diversity  Serendipity  Novelty  Relevancy

# Recommend: Hybrid Approach

# Recommend: Hybrid Approach



Hello! You are currently listening to...

ARTIST: ['SAMMY ADAMS']
ALBUM: ['BOSTONS BOY']
TRACK: ['COAST 2 COAST']

We recommend checking out...

| | artist | album | track | follower flag | review flag |
|---|---|---|---|---|---|
| 0 | SADAT X | WILD COWBOYS | THE FUNKIEST | 1.0 | 0.0 |
| 1 | N.O.R.E. | CRACK ON STEROIDS | LEHHHGOOO | 1.0 | 0.0 |
| 2 | RAH DIGGA | CLASSIC | YOU GOT IT | 1.0 | 0.0 |
| 5 | KOOL KEITH | BLACK ELVIS / LOST IN SPACE | THE GIRLS DONT LIKE THE JOB | 1.0 | 0.0 |
| 3 | CLYDE CARSON | S.T.S.A. | LET ME KNOW | 0.0 | 0.0 |
| 4 | DOE B | DEFINITION OF A TRAPPER 3 | YOU DONT EVEN KNOW | 0.0 | 0.0 |
| 6 | SAMMY ADAMS | BOSTONS BOY | COAST 2 COAST | 0.0 | 0.0 |
| 7 | PROJECT PAT | LAYIN DA SMACK DOWN | CHOOSE U | 0.0 | 0.0 |
| 8 | CASEY VEGGIES | LIFE CHANGES | YOUNG WINNERS | 0.0 | 0.0 |
| 9 | AZEALIA BANKS | FANTASEA | NEPTUNE | 0.0 | 0.0 |

# **Recommend:** Findings and Next Steps

| Findings | Next Steps |
|---|---|
| Audio, lyrics are insufficient on their own | Front End |
| User data is priority #1 | Dimensionality Reduction (Clustering, PCA) |
| Relevance + Novelty = Serendipity | User Profiles |
| High dimensionality can adversely impact engine | A/B Testing |

# Appendix

# Summary

| | Granularity | Basic Fields | Unique Fields |
|---|---|---|---|
| Spotify | Track | • Track Title<br>• Track ID<br>• Artist Name<br>• Artist ID<br>• Album Name<br>• Album ID<br>• Date | • Spotify Proprietary Audio Features<br>  • Speechiness<br>  • Liveness<br>  • Acoustiness<br>  • Danciness<br>  • Etc… |
| GENIUS Lyrics | Track | • Track Title<br>• Track ID<br>• Artist Name<br>• Artist ID<br>• Album Name<br>• Album ID | • Song Lyrics<br>• Producers |
| GENIUS User Data | Artist | • Artist ID<br>• Artist Name | • Follower ID<br>• Follower Name |
| YouTube | Album | • Artist Name<br>• Album Name | • Review Text<br>• Review Score |

| Area of Focus | Questions | Learning |
|---|---|---|
| Scraping | *Is Selenium an effective tool for bespoke scraping needs? What are the pros and cons? | Selenium is a highly customizable tool for automated testing and more difficult scraping tasks. However, it can be slow and is susceptible to interruption (ads, site scripts). |
| Merging | *What are best practices for fuzzy matching? | A number of different distance-based metrics can be leveraged for fuzzy text matching. Selecting the right distance metric depends on context. To improve match rates, leverage these features to train a classification model. |
| Hip Hop | *How has the genre changed over time? *Is hip hop "dead" as an artform? | It may be "dead" as a lyrical art form. Far less rappers are talking about their skills, focusing instead on "money" and "the hustle" in their raps. The genre has also seen a sizeable decline in rhyme density and vocab complexity. |
| NLP | *What is the most effective method for topic extraction? *What are the challenges? | Latent Dirichlect Allocation is the most popular approach here. This approach requires some subjectivity when it comes to setting the number of topics and interpreting said topics. However, there are several scores (Coherence, Perplexity) that can assist. |
| Audio | *Can we classify production based on isolated audio features? | Production is usually a team effort, which makes this task more challenging to tackle. However, for single-producer tracks, we were able to improve our baseline accuracy from 4% to ~38%. For our next iteration of this model, we'll look to leverage Mel Spectograms and a Convoolutional Neural Network. |
| Recommendations | *How do recommendations differ between collaborative filter vs content based filter? | In our implementation, content based recommendations produced more novel tracks. This is not surprising, as features like rhyme density and positive sentiment can unearth less positive tracks. |
| Recommendations | *What are the major pitfalls for recommendation engines? | *Subjectivity: The idea that content cannot capture subjective info like humor, points of view (collaborative is better here) *Scalability: as features and the user base grows, the computational cost of recommendation algorithms grows significantly. *Sparsity: In particular with user-based features this is a problem. If you are leveraging user interactions with a specific page for site recommendations, for example, it may be that only 1% of users ever get to that page at all. |
| Recommendations | *What are some methods for hybrid recommendation engines? | TBD |

| Bucket | Requirements | Advantages | Disadvantages |
| --- | --- | --- | --- |
| **Knowledge-Based** | • Significant knowledge of the domain in question | • Minimal reliance on data or specific user preferences.<br>• Approach for high complexity/barrier items | • Need extensive domain knowledge<br>• Scalability a factor |
| **Heirarchical** | • Clean product/content management system | • Simple to implement (I bought a printer, I need ink) | • Less effective for more general use cases |
| **Content Based Filtering** | • Clean and consistent metadata at the item level<br>• Optimal: user preference data | • No domain knowledge required<br>• No community required | • Overengineering possible, leading to **unexpected** results<br>• Best when aligned to user preferences |
| **Collaborative Filtering** | • Large historical data set<br>• Set of users and preferences | • No domain knowledge required<br>• Simple implementations can be effective | • Sparsity of Data<br>• Cold-Start User Problem |

| Bucket | Requirements | Advantages | Disadvantages |
|---|---|---|---|
| **Knowledge-Based** | • Significant knowledge of the domain in question | • Minimal reliance on data or specific user preferences.<br>• Approach for high complexity/barrier items | • Need extensive domain knowledge<br>• Scalability a factor |
| **Heirarchical** | • Clean product/content management system | • Simple to implement (I bought a printer, I need ink) | • Less effective for more general use cases |
| **Content Based Filtering** | • Clean and consistent metadata at the item level<br>• Optimal: user preference data | • No domain knowledge required<br>• No community required | • Overengineering possible, leading to **unexpected** results<br>• Best when aligned to user preferences |
| **Collaborative Filtering** | • Large historical data set<br>• Set of users and preferences | • No domain knowledge required<br>• Simple implementations can be effective | • Sparsity of Data<br>• Cold-Start User Problem |