Metadata problems in the music industry (and how to match music metadata)

Sertan Şentürk – Lead Data Scientist, R&D Kobalt Music Group, London

Berlin MIR Meetup Monday, September 30, 2019



- Introduction
- Behind the Scenes
- Metadata Problems
- Metadata Matching
- Summary

K Disclaimers

- I will try to cover a lot of bases
 - But we will only touch the surface

- I will not have much time to be specific for each case
 - Don't hesitate at the moment to ask for more

- I wish to make more people aware
 - So I can get lazier ^_^

K ABOUT ME

- Since 2018: Data Scientist @ Kobalt Music Group
- Previously (& briefly) @ SoundCloud
- 2017: PhD @ Music Technology Group, Universitat Pompeu Fabra
- 2011: MSc @ Georgia Tech Center for Music Technology







Behind the scenes of the music industry

K MUSIC INDUSTRY IN BRIEF (!)





K MULTI-LAYERED OPS



K OPAQUENESS







Metadata Problems

Name	Surname	Country	Genre	Role
Roger	Taylor	UK	Rock	Drums



Roger (Meddows) Taylor *Queen*



Roger (Andrew) Taylor Duran Duran

Name	Surname	Country	Genre	Role	Act	
Roger Meddows	Taylor	UK	Rock	Drums	Queen	
R.	Taylor	UK	Rock	Drums	Live Aid '85	П
Roger	Taylor	England	Rock'n Roll	Percussion	Duran Duran	I
Roger Meddows	Taylor	UK	Rock	Guitar	The Cross	Ν
Roger	Taylor	?	None	-999		Ν
Roger Meadows	Tylor	UK	Rock	Drums	Queen	E
Roger Andrew	Taylor	UK	Rock	Drums	Queen	C

INSUFFICIENT INCONSISTENT MULTI-FACETED MISSING ERRORS CONFLICTS

Name	Surname	Country	Genre	Role	Act		
Roger Meddows	Taylor	UK	Rock	Drums	Queen		
R.	Taylor	UK	Rock	Drums	Live Aid '85		
Roger	Taylor	England	Rock'n Roll	Percussion	Duran Duran	_	
Roger Meddows	Taylor	UK	Rock	Guitar	The Cross		
Roger	Taylor	?	None	-999		_	
Roger Meadows	Tylor	UK	Rock	Drums	Queen		
Roger Andrew	Taylor	UK	Rock	Drums	Queen		
HOW CAN WE FIX THIS MESS?!?							







IDEAL" SOLUTION

USE UNIVERSALLY UNIQUE IDENTIFIERS

MusicBrainz ID	Name	Surname	Country	Genre	Role	Act
558302b1-07ae-4be1-9edb-2a2b1f036313	Roger Meddows	Taylor	UK	Rock	Drums	Queen
f2b7fc94-1e8d-4f4d-9a3f-620c1a5ac564	R.	Taylor	UK	Rock	Drums	Live Aid '85
f2b7fc94-1e8d-4f4d-9a3f-620c1a5ac564	Roger	Taylor	England	Rock'n Roll	Percussion	Duran Duran
558302b1-07ae-4be1-9edb-2a2b1f036313	Roger Meddows	Taylor	UK	Rock	Guitar	The Cross
f2b7fc94-1e8d-4f4d-9a3f-620c1a5ac564	Roger	Taylor	?	None	-999	
558302b1-07ae-4be1-9edb-2a2b1f036313	Roger Meadows	Tylor	UK	Rock	Drums	Queen
	Roger Andrew	Taylor	UK	Rock	Drums	Queen

THEN, WHAT IS THE PROBLEM?!?!



THERE IS NO AGREED, UNIVERSAL DATA SOURCE FOR MUSIC... BUT... CAN'T WE CREATE ONE!?

MANY TRIED...







C DATA OPS

- Data modelling: Use the right model for right type of data
 - Tabular information (e.g. client records)?: **Relational Databases**
 - Unstructured/semi-structured data (e.g. tags)?: NoSQL Document Stores
 - Relations, communities (e.g. chain of title)?: Graph Databases
 - NOTE: There are many more things to consider! Just ask any information management person!
- Quality
- Matching
- Enrichment
- Reconciliation
- Master Data Management

K DATA OPS

- Data modelling
- Quality: Monitor the state of your metadata
 - **Completeness:** How much information is available for a document, field or dataset?
 - Coverage: How well does a data source represent others?
 - Does the data fulfil the expected types/constraints?
 - ...
 - A nice tool: <u>Amazon Deequ</u>
- Matching
- Enrichment
- Reconciliation
- Master Data Management



- Data modelling
- Quality
- Match and connect metadata as much as possible
 - Deduplicate documents (Sertan Şentürk vs. Sertan Sentuerk)
 - Link relevant documents (recordings → works)
 - Embrace linked data!
- Enrichment
- Reconciliation
- Master Data Management



- Data modelling
- Quality
- Matching
- Enrich your metadata whenever you have the chance
 - Fill missing values
 - Correct/standardize values
 - Fetch information from other sources
 - ...
 - A nice tool: <u>http://holoclean.io/</u>
- Reconciliation
- Master Data Management



- Data modelling
- Quality
- Matching
- Enrichment
- Reconciliation
 - Merge information during deduplication
 - Detect and resolve conflicting information: ontologies are a great help!
 - Ensure the processes are not creating mess
- Master Data Management



- Data modelling
- Quality
- Matching
- Enrichment
- Reconciliation
- Master Data Management
 - Keep a track of all changes, sources etc.
 - Maintain "golden" records
 - Provide a single point of truth (as much as possible)

- Cleaning everything by hand is impossible/impractical
 - Heavy monitoring and automation
- Yet keep the humans (experts with domain knowledge) in the loop:
 - to deal with difficult/edge cases
 - to evaluate your automations
- Problems will always linger
 - We have to minimize its impact
- You cannot make your data clean and shiny in a single go
 - Fix iteratively
 - Each improvement will facilitate others
- Play the long game
 - Turn it into a habit/culture of the organization

Κ **ULTIMATE GOAL: OPEN & LINKED DATA**



METADATA MATCHING

K METADATA MATCHING

- Has many names in the literature
 - record linkage, data integration, entity recognition, identity resolution, ...
- People mean different things
 - e.g. joining tables, enrichment, conflict resolution, ...

 Metadata matching: Creating a (named) link between documents which are represent the same entity (same person in two different datasets) or have a relationship (recording to its work)

- Retrieval problem
 - Sometimes ranked
- Targets vs Queries
- (Typically) binary relevance
- Cardinalities vary depending on the task
 - Recording to work: low
 - Work to recording: higher

EVALUATION

- Depending on the cardinality and use case, e.g.
 - Human (MAP, MRR) vs. Fully automatic (F₁); needle in haystack (Prefer PR over ROC-AUC); Low tolerance to false positives (Weight precision more than recall)



K EVERYONE HAS SIMILAR PROBLEMS

- Insurance
- Finance
- Manufacturing
- Medicine

. . .

Online stores

hm ok thanks google Jan van Balen Painter PERIOD PEOPLE ALSO SEARCH FOR OVERVIEW Jan van Balen was a Flemish painter known for his Baroque paintings of history and allegorical subjects. He also painted landscapes and genre scenes. Wikipedia Born: 21 July 1611, Antwerp, Belgium Died: 14 March 1654, Antwerp, Belgium Period: Baroque Parents: Hendrick van Balen

Follow

 \sim

10:21 am - 28 Sep 2019

Jan Van Balen

@jvanbalen

Usage processing: Match <u>consumption</u> reported by streaming services to <u>works</u>

Royalty processing: Match royalty statements to works

Onboarding: Merge entities (<u>creators, works, recordings</u>) in an external catalogue with the internal catalogue

Deduplication: Find <u>duplicate</u> entities inside a catalogue

Enrichment: Match entities in an external catalogue with the internal catalogue to <u>complete</u> <u>missing information</u>

Relation discovery: Identify <u>missing links</u> between entities (e.g. collaborators in an album)

Error resolution: ...

SOLUTIONS

RULE-BASED SOLUTIONS

- First choice to go
 - Before big data era (till late 2000s)
 - (many music businesses are still here)
- Based on simple, deterministic decisions
 - e.g. full/sub string match between separate fields and converting the results to a "match percentage" according to some business logic
- Seems interpretable at first, but rules/logic explode as time passes
 - Spurious and/or conflicting rules
- Does not scale beyond small data (>10k documents)
- Does not generalize to different tasks
- Don't worry if you are at this stage
 - You can use the system to <u>collect initial data</u> to **evaluate** future methods

COUT OF BOX SOLUTIONS

- If you don't have engineering resources or time
- No need to build or maintain much infrastructure



- <u>"FindMatches" ML transform</u>
- Works hand-to-hand with Lake Formation
- Pioneer DataOPS start-up
- Does more than matching







Şentürk, S. (2019). Music metadata matching using ElasticSearch. London ElasticSearch Meetup

K FILTERING USING ML



Nice tool: <u>Snorkel</u> for weakly labelling data

WHERE IS DEEP LEARNING!?

- Late adoption in the academia & industry
- "deepmatcher" (Sigmod 2018)
 - How does DL architectures fare against SOTA?
 (Does not outperform for structured data yet)
 - SIF, RNN, attention-based & hybrid models using word embeddings
 - Code (<u>Github</u>)
 - Also evaluated on (rather small-sized) music metadata (<u>Demo</u>)



Mudgal, Sidharth, et al. "Deep learning for entity matching: A design space exploration." *Proceedings of the 2018 International Conference on Management of Data.* ACM, 2018.

GRAPH-BASED APPROACHES



- node2vec: <u>http://snap.stanford.edu/node2vec/</u>
- Graph Convolution Neural Networks

Oramas, S., & Sordo, M. (2015) Knowledge acquisition from music digital libraries. IAML/IMS Congress



- Music metadata is (to put it kindly) bad...
- Why?!?
 - uncoordinated supply chains, data quality issues, absence of trust, lack of data literacy, human errors, and national politics...
- What should we do?
 - First, clean your own turf
 - Follow information management best practices and embrace linked data principles
 - Connect your metadata with others

K SUMMARY (2)

- Metadata matching is inevitable in the current ecosystem (and till eternity)!
 - i.e. symptom of lack of interconnectivity between metadata resources
 - If you sweep your data problems under the carpet, they will hunt you later in the form of matching...
 - Remedy your problems before they become a sickness requiring a surgery...
- Where should I start?
 - Focus on a single use case; generalize later
 - Spend time on understanding the use case, pick a relevant eval measure, collect data
 - Blocking using ElasticSearch; use simple ML models (e.g. tree-based) based on TF-IDF to filter candidates



- Conferences: <u>Sigmod</u>, <u>VLDB</u>, <u>SigKDD</u>
- Academics: Michael Stonebraker (MIT CSail), Ihab Francis Ilyas (Uni Waterloo), AnHai Doan

& <u>Theodoros Rekatsinas</u> (UW Madison), <u>Xu Chu</u> (Georgia Tech), <u>Christopher Ré</u> (Stanford)

- Course Material: <u>Duke Uni CompSci590.01</u>, <u>UW Madison CS 838</u>
- Companies: <u>Tamr</u>, <u>Trifacta</u>
- Tools: AWS Lake Formation, ElasticSearch, Apache TinkerPop

Magellan, deepmatcher, py stringmatching, py stringsimjoin

Snorkel, HoloClean

<u>Deequ</u>

FURTHER READING

- Stonebraker, M., & Ilyas, I. F. (2018). <u>Data integration: The current status and the way</u> <u>forward</u>. IEEE Data Eng. Bull., 41(2), 3-9.
- Hellerstein, J. M. (2008). <u>Quantitative data cleaning for large databases</u>. United Nations Economic Commission for Europe (UNECE).
- Oramas, S., & Sordo, M. (2015) <u>Knowledge acquisition from music digital libraries</u>.
 IAML/IMS Congress
- Allemang, D., & Hendler, J. (2011). <u>Semantic web for the working ontologist: effective</u> modeling in RDFS and OWL. Elsevier.
- Kipf, T. (2016). <u>Graph convolutional networks</u>. Self-Published
- Şentürk, S. (2019). <u>Music metadata matching using ElasticSearch</u>. London ElasticSearch Meetup

We are hiring!!!

Lots of positions and roles, incl. data engineers at all levels

Speak with me

https://www.kobaltmusic.com/ company/careers



