# A Semantic Landscape of the Last.fm Music Folksonomy

**Joseph Biberstine, Russell J. Duhon, Elisha Allgood, Katy Börner**
*Cyberinfrastructure for Network Science Center,*
*School of Library and Information Science,*
*Indiana University*
**André Skupin**
Department of Geography,
*San Diego State University*

INDIANA UNIVERSITY

cyberinfrastructure for NETWORK SCIENCE CENTER
cns.iu.edu

SAN DIEGO STATE UNIVERSITY

# Motivation

- Domain
  - How is the world of music and music experience organized?
    - What kinds of themes emerge in this domain and what is their structure?

- Challenges
  - Collect and prepare high-dimensional social data
  - Create a model large enough to faithfully represent the domain
  - Train a model of this substantial size
  - Design a visualization that does justice to the richness of the model

folk · post-rock · post rock · new age · idm · techno · ambient · electronic · electronica · dark ambient · drone · experimental · experimental metal · neofolk · chillout · singer-songwriter · indie · indie pop · shoegaze · krautrock · avant-garde · electro · drum and bass · trip-hop · female vocalists · alternative · indie rock · deathrock · gothic rock · darkwave · industrial · ebm · industrial rock · dark electro · female vocalists · french · soundtrack · alternative rock · britpop · british · post-punk · new wave · industrial metal · german · japanese · classical · piano · punk · punk rock · hardcore · ska · new wave · disco · pop · dance · latin · world · 80s · classic rock · glam rock · hard rock · post-hardcore · screamo · christian · cover covers · italian · 90s · soul · funk · psychedelic · guitar virtuoso · 60s · 70s · heavy metal · nu metal · metalcore · grindcore · country · rap · hip-hop · jazz · swing · progressive rock · progressive · progressive metal · power metal · thrash metal · heavy metal · nwobhm · folk metal · symphonic metal · sludge · death metal · melodic death metal · black metal · reggae · blues · guitar · doom metal · gothic metal

# Raw Data - Source

- Last.fm is a social Internet radio site
  - Users share information about songs they are listening to
  - They can also tag songs
    - With any strings of text they like

**last.fm**

**Need new music?**

Last.fm lets you effortlessly **keep a record of what you listen to*** from any player. Based on your taste, Last.fm recommends you more music and concerts!

*We had to invent a word for this, it's called scrobbling.

# Raw Data - Summary

- Gathered during the first half of 2009
- 99,405 registered users
  - 52,452 active
- 281,818 tags
- 1,393,559 songs

- 10,936,545 annotations
  - An annotation is a (user, tag, song) triple, a tagging event

Data originally collected for:
Schifanella, R., Barrat, A., Cattuto, C., Markines, B., and Menczer, F. (2010). Folks in Folksonomies: Social Link Prediction from Shared Metadata. Proc. 3rd ACM International Conference on Web Search and Data Mining (WSDM).

# Top Tags

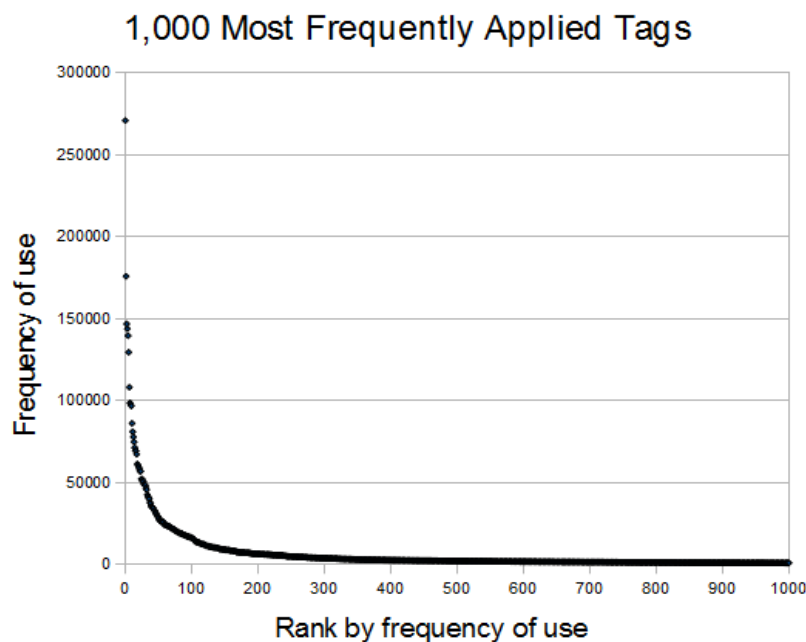| | | |
|---|---|---|
| rock | singer-songwriter | heavy metal |
| electronic | 80s | chillout |
| seen live | folk | dance |
| indie | hard rock | british |
| alternative | progressive rock | 90s |
| pop | indie rock | psychedelic |
| female vocalists | electronica | blues |
| jazz | punk | hip-hop |
| classic rock | instrumental | post-rock |
| experimental | soul | new wave |
| ambient | black metal | soundtrack |
| metal | industrial | classical |
| alternative rock | death metal | 00s |

# Tags Are More Than Just Genres

- Intensional
  - From recognized genres to simple objective facts
    - *rock* (rank 1)
    - *electronic* (2)
    - ..
    - *female vocalists* (7)
      - *female vocalist* (64)
    - *acoustic* (51)
    - ..
    - *title is a full sentence* (101)
- Extensional
  - A mix of social signals, properties of the user-song experience, and aides to personal categorization
    - *seen live* (3)
    - *beautiful* (48)
    - *favorites* (54)
    - *albums i own* (97)
    - *altar of the metal gods* (58)
      - A case of graffiti?

# Raw Data - Thresholding

- The self-organizing map (SOM) method will not scale to 280,000+ tags/dimensions in raw form
  - Not often used with more than hundreds of dimensions
- Consider only the 1,000 most frequently applied tags
  - Keep only songs annotated by some user with any of these tags



1,000 Most Frequently Applied Tags

# Thresholded Data - Summary

|  | **Raw** | **Thresholded** |
|---|---|---|
| Tags | 281,818 | 1,000 |
| Songs | 1,393,559 | 1,088,761 (78% of original) |
| Annotations per song (average) | 7.8 | 6.8 |

# Approach

- Characterize each song as a vector over each tag dimension
  - Each coordinate is the number of annotations
    - Summed across users

- A song is a piece of tag relationship evidence

# Method - Background

- Self-organizing maps
  - Neural network training algorithm
  - Unsupervised
  - High-dimensional data

⇩

Low-dimensional discrete geometric model

  - Goal:
    - Proximity in the input space

⇩

Proximity on the map

# Self-Organizing Map Algorithm - Classical

1. Create a lattice of neurons
2. To each neuron assign an initial (often random) vector with as many dimensions as the training data
3. For each training vector:
    1. Identify the neuron of minimal distance according to the input space metric (the "best-matching unit")
    2. For each neuron:
        1. Pull this neuron's vector toward the training vector in proportion with this neuron's distance from the best-matching unit

# Self-Organizing Map Algorithm - Parallelized Implementation

- A previous project trained on twice as many data and twice as many dimensions
  - Completely intractable using widely available software
  - Created our own implementation
    - Divide the training data among multiple processes
      - Each process holds a complete copy of the map
      - Periodically synchronize process-local copies of the map to create a new process-global map

- Adapted with several project-specific optimizations from:
  - Lawrence, R.D., Almasi, G.S., Rushmeier, H.E. (1999). A Scalable Parallel Algorithm for Self-Organizing Maps with Applications to Sparse Data Mining Problems. Data Mining and Knowledge Discovery

# Training the Map

- 2D hexagonal lattice of neurons
  - 180 on either side = 32,400 altogether

- Input space metric: **cosine similarity**
  - Induced interpretation: Each training vector (and so consequently each neuron vector) represents a direction in the 1,000-dimensional tag space
  -

$$\text{similarity} = \cos(\theta) = \frac{A \cdot B}{\|A\|\|B\|} = \frac{\sum_{i=1}^{n} A_i \times B_i}{\sqrt{\sum_{i=1}^{n} (A_i)^2} \times \sqrt{\sum_{i=1}^{n} (B_i)^2}}$$

- 50 complete passes over the training data

# Computation

- 300 processes across 100 compute nodes of Big Red, a supercomputer at Indiana University
- Parallel runtime = 13 hours
    - Serial equivalent runtime = **5 months**

# Legend

# Visualization

- ☐Recall there is a corresponding vector to each neuron which describes its position in the input space
- In other words, its position along each tag dimension

- Consider the $n^{th}$ strongest tag association of each neuron
- A contiguous swath of neurons sharing a common $n^{th}$ strongest tag association is termed a region
- As the map is trained over 1,000 tags, we have 1,000 distinct partitions of the map into such regions

folk · post rock · post-rock · new age · idm · techno · ambient · experimental · dark ambient · drone · electronic · indie · electronica · singer-songwriter · indie pop · krautrock · neofolk · ebm · trip-hop · female vocalists · alternative · female vocalist · cabaret · shoegaze · avant-garde · noise · experimental metal · industrial · chillout · indie rock · deathrock · loved · darkwave · ebm · industrial rock · industrial metal · dark electro · female vocalists · french · alternative rock · hard rock · britpop · british · gothic rock · post-punk · german · japanese · soundtrack · alternative · alternative rock · alternative · punk · new wave · female vocalists · latin · classical · rock · britpop · punk rock · new wave · disco · world · piano · dance · instrumental · 80s · hardcore · ska · pop · 90s · classic rock · glam rock · hard rock · post-hardcore · screamo · christian · cover · italian · soul · funk · psychedelic · 60s · classic rock · 70s · hard rock · heavy metal · nu metal · metalcore · grindcore · hip-hop · jazz · guitar · southern rock · progressive rock · blues rock · hard rock · heavy metal · nwobhm · folk metal · metal · country · rap · hip hop · swing · blues · guitar · progressive rock · progressive · progressive metal · power metal · thrash metal · symphonic metal · sludge doom · death metal · black metal · reggae · progressive · power metal · gothic metal · doom metal · melodic death metal · altar of the metal gods

# Interpretation

- Interpreters report a mix of
  - Recognition
    - Patterns of hierarchical and neighborhood relationships among tags match expectations
  - Discovery
    - Opportunities to find new musical categories
  - Surprise
    - Relationship between *rock*, *blues*, and *jazz*

darkwave industrial rock dark electro female vocalists french

ethereal goth darkwave goth ambient synthpop ebm metal heavy metal ebm industrial dark electro alternative female pop vocalists

ethereal electronic synthpop industrial metal german female vocalists brazilian bossa nova

new wave synthpop synthpop deutsch female vocalists rock pop female vocalist latin

new wave synthpop synthpop gay disco dance soul female vocalists spanish arabic

post-punk new wave synthpop soul disco funk dance dance disco bossa nova latin salsa spanish

classic rock pop soul disco 90s dance electropop female vocalist soul

post-hardcore emocore pop rock britpop comedy 90s 90s rnb

screamo emo all the best alternative singer-songwriter male vocalists urban

post-hardcore gospel christian covers italian folk male vocalist dance

screamo christian italian italiana rnb chill r&b

rap hip hop eurovision fun 2008 france underground hip-hop swedish hip hop french italian urban country rap new york old school hiphop russian polish male vocalists female vocalists female vocalist underground hip-hop garage spoken word funk i am a party girl here is my soundtrack greek ninja tune turntablism r&b old school smooth jazz jazzy british rap comedy known yet able artists polish pop folk nu jazz smooth hiphop hip-hop finnish spanish metallic gothic metal norwegian underground hip-hop rap dancehall rap hip hop sad urban oldies love memories 70s 80s 00s 90s 2008 dance political new school usa funky

industrial rock · industrial · ebm · trip-hop · trip hop · downtempo · soundtrack

industrial rock · dark electro · gothic · female vocalists · chillout · japanese · j-rock

synthpop · industrial metal · dark electro · cabaret · french · chanson francaise · beautiful · instrumental

german · deutsch · female vocalists · canadian · bossa nova · mpb · brasil · piano · classical

disco · funk · female vocalists · british · latin · brazilian · world · african · romantic · violin

soul · dance · spanish · arabic · africa · love · opera · classical · baroque

80s · 90s · pop · electropop · sexy · amor · northern soul · motown · rhythm and blues · rnb · soul

alternative · singer-songwriter · britpop · folk · dance · rnb

# Potential Applications

- Interactive music navigator and playlist generator
- Mapping portfolios as fields of neuronal activation
  - For the set of songs associated with any entity, we can see where in the world they belong
    - A user: Their favorite songs
    - A band: Their complete work
    - A group of users: What is their turf?
  - .. or look at the difference of any of these fields
    - What is the difference between The Who and The Guess Who?
    - How has this entity moved through the world of music over time?
    - Where have listeners like me headed next?

# Contributors

- Joseph Biberstine, Indiana University
    - Contact: jrbibers@indiana.edu
- André Skupin, San Diego State University
- Russell J. Duhon, IU
- Elisha Allgood, IU
- Katy Börner, IU