

# Using TDM Studios to help systematically analyze research in music

Colloquium on Text & Data Mining in Libraries 2023

# Welcome to TDM Studio



## Workbench Dashboard



### WHAT TO EXPECT

- Requires basic coding in Python or R and text analysis
- Provides options to select content and create datasets
- Offers Jupyter Notebook coding platform

#### Predict Labels

By using SBERT model and the array of parsed sentences, we will now output a corresponding array where each element is a tuple of the predicted label and an array of the raw possibilities for each label.

Sample code for text analysis

```
# Set sbert_path to location of SBERT model
sbert_path = 'Overall_emotion_classifier/nli-mpnet-base-v2'
transformer = ST(sbert_path)

# Encode the parsed_sents array
sentence_embeddings = transformer.encode(parsed_sents, show_progress_bar=False)

# Both scaler and sentiment_model should exist before running this cell
if scaler is not None and sentiment_model is not None:
    standardized = scaler.transform(sentence_embeddings)

    y_pred_numeric = sentiment_model.predict(standardized)
    y_pred_string = label_encoder.inverse_transform(y_pred_numeric)

    # Call the predict function on our sentences
    raw_predictions = sentiment_model.predict_proba(standardized)

    results = list(zip(y_pred_string, raw_predictions))

    # Print first element of array
    print(results[0:5])
else:
    print("Please load scaler and sentiment model.")
```

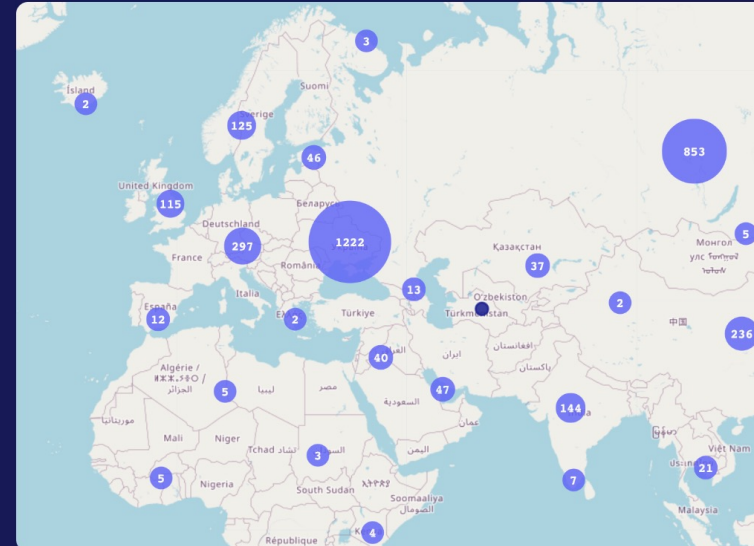


## Visualization Dashboard



### WHAT TO EXPECT

- Requires no coding experience
- Provides the ability to select content specific to your project
- Offers pre-configured data visualization to perform analysis



a187 UNIVERSITY OF TORONTO

Open Jupyter Notebook



Running



Toggle to restart your environment

## Datasets (Using 1 of 10)

Create datasets of up to 2,000,000 documents per dataset. When a dataset is completed, it will transfer automatically to the folder on your Jupyter Notebook.

<input type="checkbox"/>	NAME	DATE RANGE	SEARCH ?	DATA SOURCE ?	DOC COUNT
<input type="checkbox"/>	dissertations93-2002	Jan 01, 1993 to Dec 31, 2020		ProQuest Dissertations & Theses Global	891,443

+ Create New Dataset

Publication Titles

ProQuest Databases

Congressional Hearings **BETA**



## &lt; Choose Databases (313)

dissertations

x



\* No results found for "dissertations"

<input type="checkbox"/>	ABI/INFORM Collection	The most comprehensive ABI/INFORM™ database, this comprises ABI/INFORM Global, ABI/INFORM Trade and Industry, and ABI/INFORM Dateline. The database features thousands of full-text journals, dissertati...	
<input type="checkbox"/>	Canadian Business & Current Affairs Database	This database is a longstanding, comprehensive Canadian periodical collection covering multiple subjects and topics, with millions of full-text records. Accessible to readers and researchers at every ...	
<input type="checkbox"/>	Dissertations & Theses @ University of Toronto	This database gives access to the dissertations and theses produced by students at your institution.	Full text
<input type="checkbox"/>	ERIC	This database is sponsored by the U.S. Department of Education to provide extensive access to education-related literature. ERIC provides coverage of journal articles, conferences, meetings, governmen...	
<input type="checkbox"/>	Linguistics and Language Behavior Abstracts (LLBA)	This database abstracts and indexes the international literature in linguistics and related disciplines in the language sciences. The database covers all aspects of the study of language including pho...	
<input checked="" type="checkbox"/>	ProQuest Dissertations & Theses Global	ProQuest Dissertations & Theses (PQDT) Global is the world's most comprehensive collection of dissertations and theses from around the world, offering millions of works from thousands of universities....	Full text
<input type="checkbox"/>	Social Services Abstracts	This database provides bibliographic coverage of current research focused on social work, human services and related areas, including social welfare, social policy and community development. The datab...	
<input type="checkbox"/>	Sociological Abstracts	Sociological Abstracts, and its companion file Social Services Abstracts, cover the international literature of sociology, social work, and related disciplines in the social and behavioral sciences. I...	
<input type="checkbox"/>	Sociological Abstracts	Sociological Abstracts, and its companion file Social Services Abstracts, cover the international literature of sociology, social work, and related disciplines in the social and behavioral sciences. I...	

1 Database selected

Next: Refine Content



The screenshot displays the JupyterLab web interface. At the top, the header shows the Jupyter logo, the text "jupyter Mus\_diss", and "Last Checkpoint: 11/26/2022 (autosaved)". A "Logout" button is in the top right. Below the header is a menu bar with "File", "Edit", "View", "Insert", "Cell", "Kernel", "Widgets", and "Help". To the right of the menu bar is a "Not Trusted" status indicator and a "conda\_python3" environment selector. Below the menu bar is a toolbar with icons for saving, adding, deleting, and running code, along with a "Code" dropdown and an "nbdiff" button. The main area is divided into three panes. The left pane shows a code editor with two input cells. The first cell, labeled "In [26]:", contains the following Python code:

```
import os
import pandas as pd
import re
import sys
```

The second cell, labeled "In [27]:", contains the following Python code:

```
19#_20_2
#_18_20
base =
input_file =
_2002 =
```

The right pane shows a file browser. It has tabs for "Files", "Running", and "Clusters". Below the tabs is a message "Select items to perform actions on them." and a search bar. The file browser displays a directory structure with a folder icon and the text "0" and "output\_files". Below this, there is a folder icon and the text "seconds ago". At the bottom, there is a file icon and the text "Untitled.ipynb", followed by the text "Running 5 months ago" and "6.9 kB".

Upload

New ▾

↻

↓	Last Modified	File size
	seconds ago	
	2 months ago	
	5 months ago	
unning	5 months ago	72 B

```

for i in input_files:
    elementList = []
    with open(base+"/"+i, "r") as f:
        lines = f.read()
        ClassCode = re.findall(r'<ClassCode>([0-9]{4})',lines)
        #metaElement(ClassCode)
        if ClassCode and "0413" in ClassCode:
            print(i)
            tags = re.findall(r'<([\w\W]*)>',lines)
            for t in tags:
                print(t)

```

Filtered out dissertations not classified as music

```

DegreeDate = re.search(r'<DegreeDate>([\w\W]*)<',lines)
metaElement(elementList,DegreeDate.group(1))
RawLang = re.findall(r'<RawLang>([\w\W]*)<',lines)
metaElement(elementList,RawLang)
Title = re.findall(r'<Title>([\w\W]*)<',lines)
metaElement(elementList,Title[0])
ClassExpansion = re.findall(r'<ClassExpansion>([\w\W]*)<',lines)
metaElement(elementList,ClassExpansion)
SchoolLocation = re.findall(r'<SchoolCodeName>([\w\W]*)<',lines)
metaElement(elementList,SchoolLocation[0])
GenSubjTerm = re.findall(r'<GenSubjValue>([\w\W]*)<',lines)
metaElement(elementList,GenSubjTerm)
DissPaperCategory = re.findall(r'<DissPaperCategory>([\w\W]*)<',lines)
metaElement(elementList,DissPaperCategory)
DissPaperKwd = re.findall(r'<DissPaperCategory>([\w\W]*)<',lines)
metaElement(elementList,DissPaperKwd)
GenSubjValue = re.findall(r'<GenSubjValue>([\w\W]*)<',lines)
metaElement(elementList,GenSubjValue)

```

Metadata elements were gathered by their XML tags

- Title
- Degree date
- Category





- Dataset Export data as csv – to excel

degree	language	title	Classification	institution
1995	['English', 'English']	Vocal lyricism in the melodies of Andre Caplet: Two lecture recitals revealing a singer's perspective	['Music', 'Music education']	University of Maryland, College Park
1995	['English', 'English']	Vocal fatigue in choral singing: Causes and suggestions for prevention voiced by prominent choral directors	['Music', 'Music education', 'Speech therapy']	The Florida State University
1995	['English', 'English']	Violin concerto: For violin and computer-generated tape. (Original composition);	['Music', 'Music education']	Stanford University
1995	['French', 'French']	Une approche poietique du volet composition en education musicale	['Music education', 'Music', 'Cognitive psychology']	Universite de Montreal (Canada)
1995	['English', 'English']	Twentieth century music for unaccompanied trumpet: An annotated bibliography	['Music', 'Music education']	Louisiana State University and Agricultural & Mechani
1995	['English', 'English']	Twentieth century Chinese vocal music with particular reference to its development and nationalistic characterist	['Music education', 'Music']	University of California, Los Angeles
1995	['English', 'English']	Tracking the nature of melodic expectancy development in musical children	['Music', 'Developmental psychology', 'Music education']	University of Washington
1995	['English', 'English']	The value of analysis in the maturation of cognitive musicianship: An experiential chronicle	['Music education', 'Music']	Duquesne University
1995	['English', 'English']	The Twelve Bagatelles and Sonata-Fantasia of George Rochberg: A performer's analysis	['Music', 'Music education', 'American studies']	Southwestern Baptist Theological Seminary
1995	['English', 'English']	The solo pianist: A critical analysis of concepts of musical giftedness	['Music education', 'Music', 'Psychology']	Concordia University (Canada)
1995	['English', 'English']	The relationships among choral performance quality, choral student emotive and aesthetic perception, and audie	['Music education', 'Music']	The University of Utah
1995	['English', 'English']	The relationship between text and music in the choral works of Robert H. Young	['Music', 'Biographies', 'Music education']	The Southern Baptist Theological Seminary

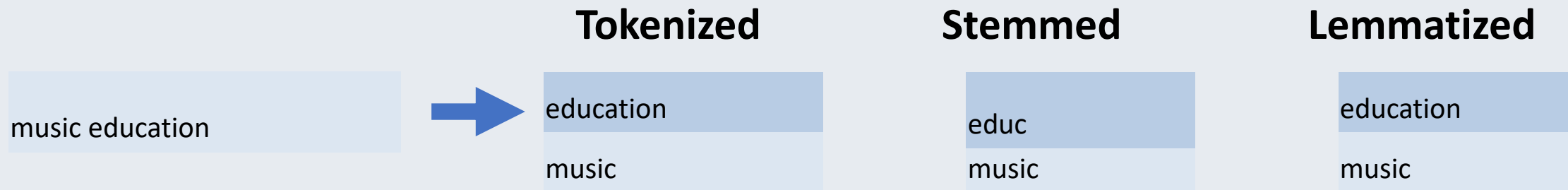
- First, classification terms were considered
- A Counter object was created for all the terms

	A	B	C	D	E
1	Dissertation year	Number of dissertations	number of classification terms	number of unique terms	terms as percent
2	2022	103	431	71	4.184466019
3	2021	165	715	90	4.333333333
4	2020	135	548	78	4.059259259
5	2019	152	612	93	4.026315789
6	2018	151	518	67	3.430463576
7	2017	161	492	70	3.055900621
8	2016	172	519	68	3.01744186
9	2015	164	521	65	3.176829268
10	2014	208	657	79	3.158653846
11	2013	187	567	65	3.032085561



Here the classification terms were :

- Normalized
- Tokenized
- Stemmed
- Lemmatized
- Counted



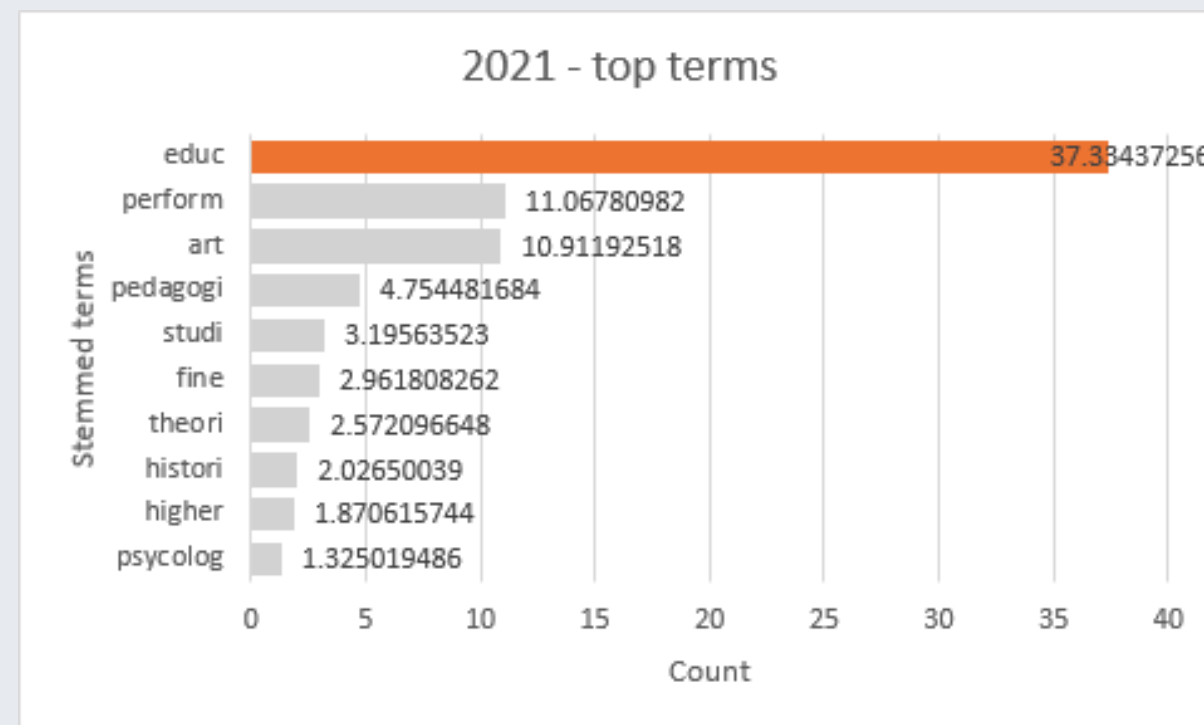
Classification Term	count	Tokenized	count	Stemmed	count	Lemmatized	count
teacher education	12	education	451	educ	479	education	451
music education	268	music	323	music	378	music	323
musical performances	43	pedagogy	61	perform	142	pedagogy	61
dance	2	performing	99	art	140	performing	99
performing arts	39	technology	8	pedagogi	61	technology	8
cultural anthropology	3	middle	4	studi	41	measurement	1
fine arts	38	musical	55	fine	38	musical	55
music history	21	development	8	theori	33	curriculum	8
reading instruction	1	gender	5	histori	26	psychology	17
latin american studies	3	curriculum	8	higher	24	anthropology	3
secondary education	9	theory	33	psycholog	17	higher	24
educational technology	8	teacher	12	teacher	12	teacher	12





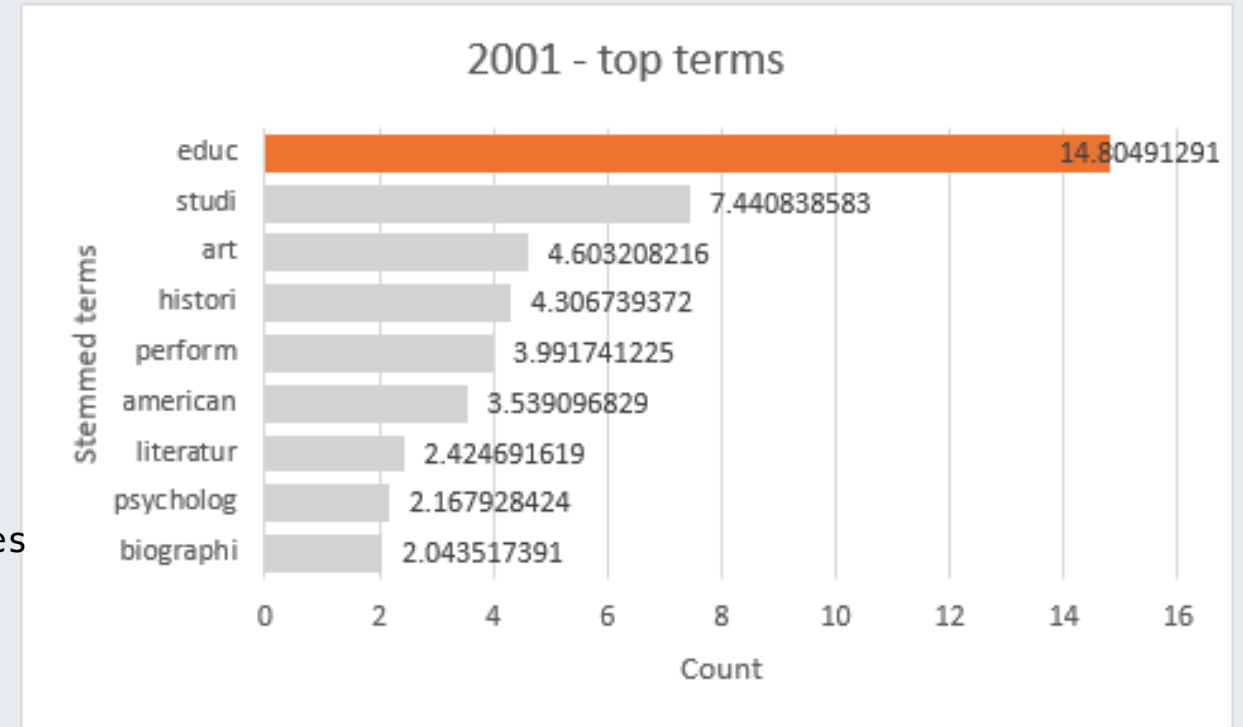
## Terms collocated under “studi” -2021

- african american studies
- asian studies
- black studies
- caribbean studies
- east european studies
- ethnic studies
- european studies
- latin american studies
- middle eastern studies
- near eastern studies
- scandinavian studies
- slavic studies
- web studies
- womens studies



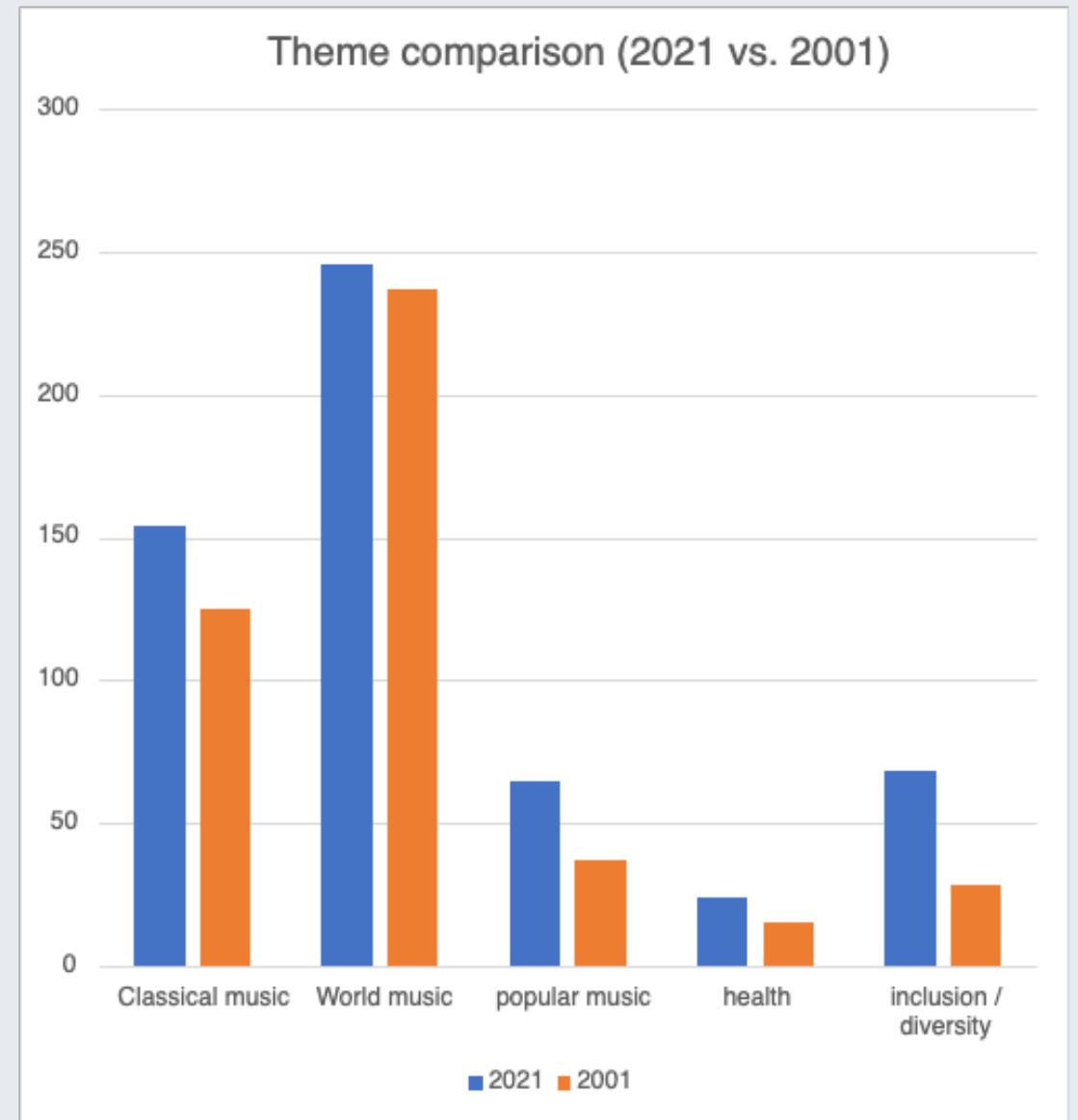
## Terms collocated under “studi” -2001

- baltic studies
- southeast asian studies
- environmental studies
- indigenous studies
- north african studies
- pacific rim studies
- scandinavian studies
- peace studies
- sound recording studios
- social studies education
- biblical studies
- disability studies
- asian american studies
- classical studies
- slavic studies
- sub saharan africa studies
- military studies
- translation studies
- east european studies
- holocaust studies
- african studies
- islamic studies
- museum studies
- hispanic american studies
- south african studies
- judaic studies
- caribbean studies
- lgbtq studies
- web studies
- canadian studies
- native american studies
- south asian studies
- regional studies
- european studies
- individual & family studies
- near eastern studies
- asian studies
- latin american studies
- middle eastern studies
- gender studies
- african american studies
- ethnic studies
- film studies
- black studies
- womens studies
- american studies

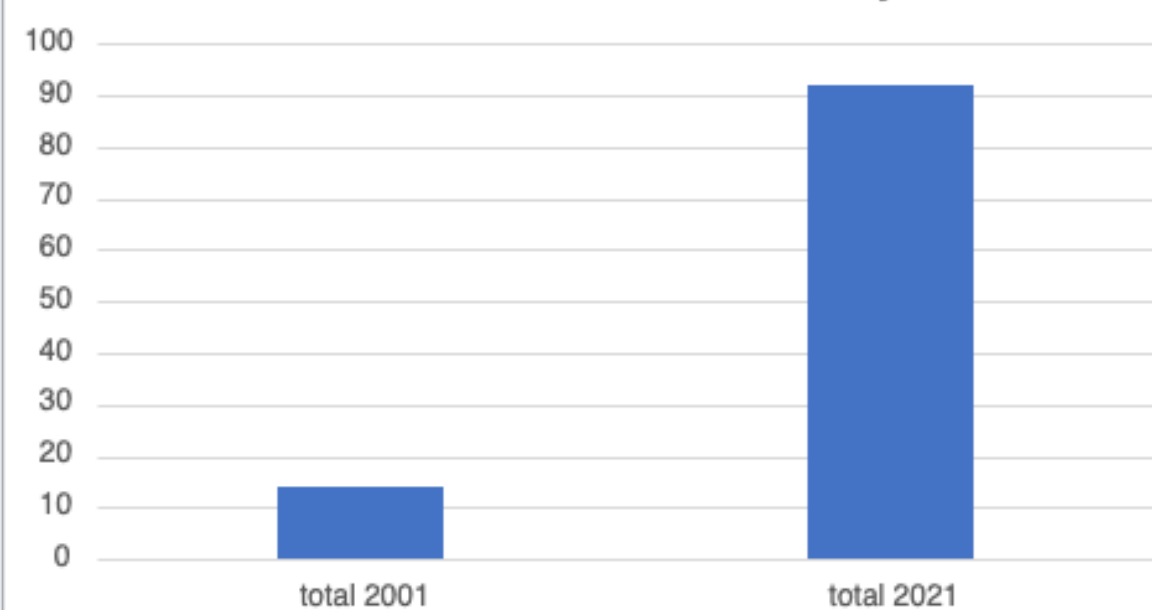


## Titles

- Tokenized
- Lemmatized
- Bigrams created
- Topic modelling



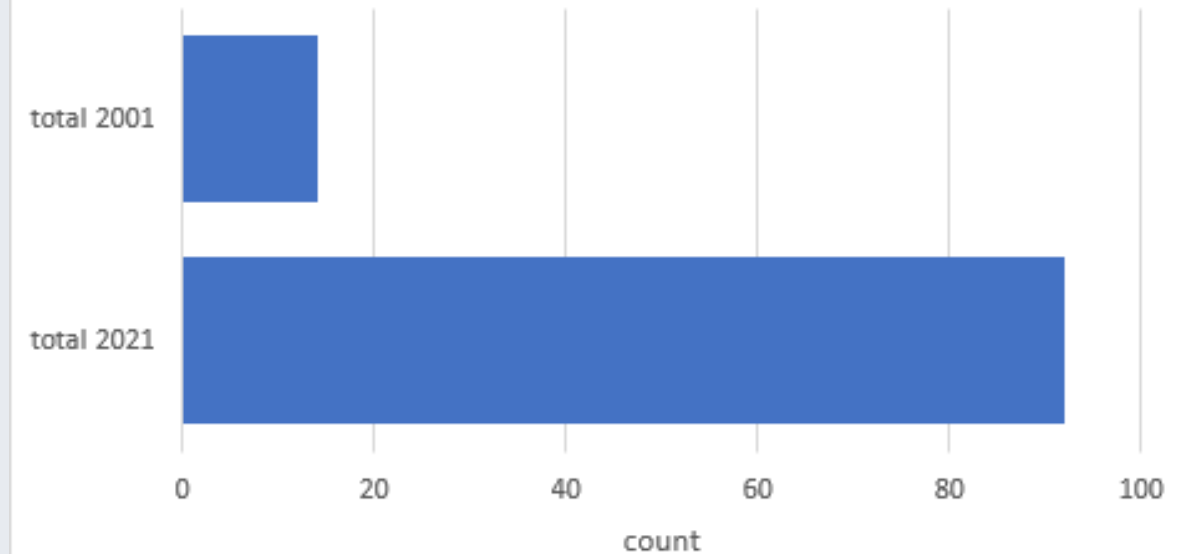
## mental health and diversity



### Mental health term examples

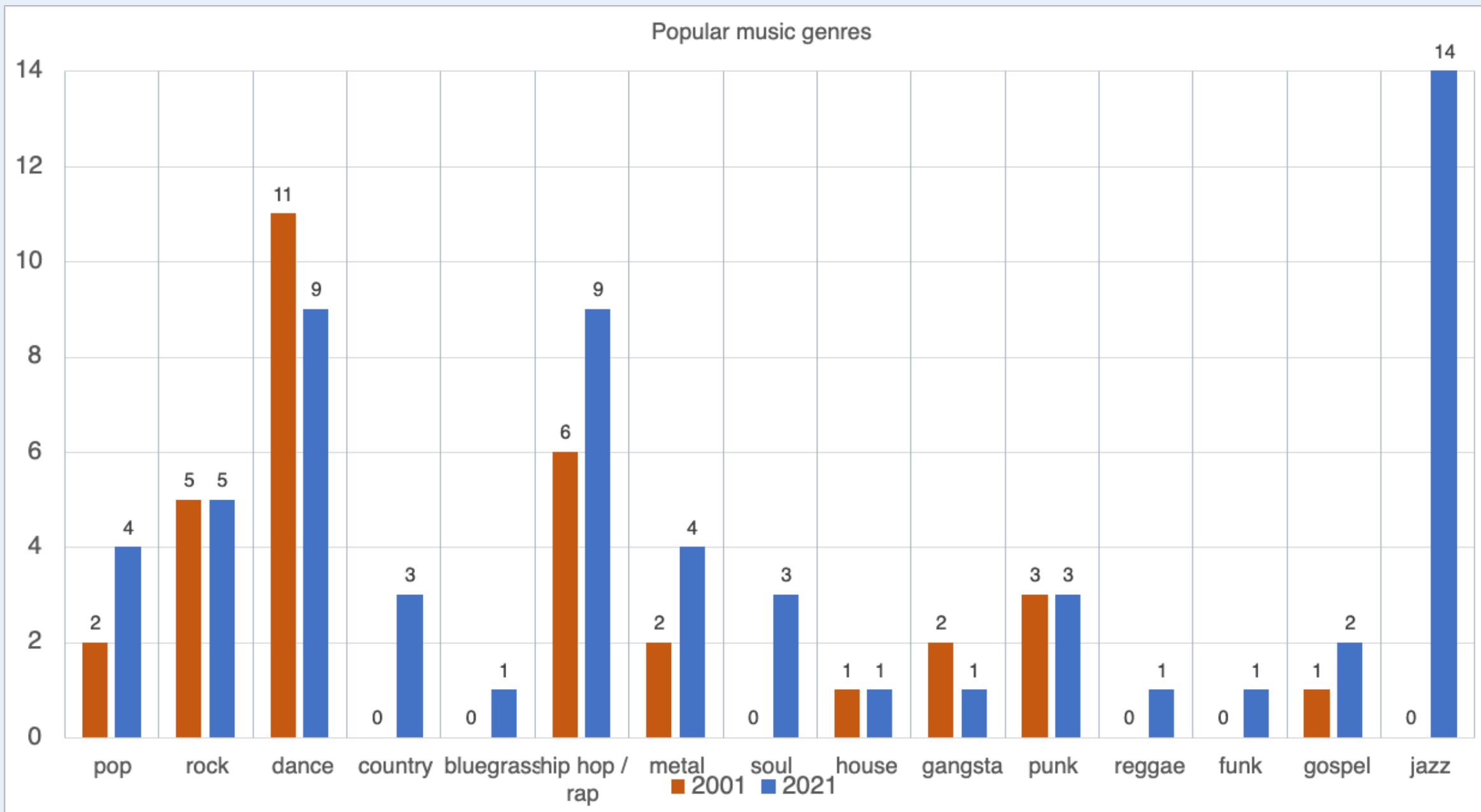
therapy  
mental health  
disorder  
depression  
crisis  
anxiety

## Popular music terms



### Inclusion term examples

justice  
identity  
bias  
intersectionality  
empathy  
gender  
feminism



# Nonwestern/art music 2021



# Nonwestern/art 2001





# Collection assessment Collection development Instruction



# Thank You

James Mason, Metadata and Digital Initiatives Librarian  
j.mason@utoronto.ca