# Using Machine Learning to Identify Non-Inclusive Language in the Library Catalog

Hollie Gardner & Ryan Kinney

Advisor: Nedelina Teneva, PhD

Hollie Gardner
MSDS Dec '22

Ryan Kinney
MSDS May '23

Dr. Nedelina Teneva
Advisor

# Team Members

# Content Disclosure

This project addresses harmful language that some audience members may find difficult to discuss. Language in this presentation has not been fully censored.

# Microaggressions

- "Everyday, subtle, intentional – *and oftentimes unintentional* – interactions or behaviors that communicate some sort of bias toward historically marginalized groups"[2]

- "Death by a Thousand Cuts"[1]

# Microaggressions

**Impact:**

- Lower GPA or researcher productivity[4]

- Reduction in retention and graduation rates[4]

- Negative impacts to mental and physical health[4,5,6]

# In the Catalog

- Presence of harmful and non-inclusive language

- Former Student: Interruption to productivity

# In the Catalog

Where did it come from?

- Historical language usage
- Rules of cataloging
- Globally shared records

## Problem:

Libraries **should be inclusive** to all groups and **should not create harm** to students and researchers from historically marginalized groups and this isn't always happening.

## Solution:

Libraries need **a proactive way to identify and remove harmful and non-inclusive language** from their millions of catalog records.

# Data: Library Catalog

# Data: Bibliographic Record, Public View

# Data: Bibliographic Record, Source View



```
leader  02026cam a22004453i 4500
001     9917689040300041
005     20220703201620.0
006     m     o   d |
007     cr cnu||||||||
008     220703s2018    xx       o     ||||0 eng d
020     ##$a9781978504431 $q(electronic bk.)
020     ##$z9781978504660
035     ##$a(MiAaPQ)EBC5733747
035     ##$a(Au-PeEL)EBL5733747
035     ##$a(OCoLC)1083356992
040     ##$aMiAaPQ $beng $erda $epn $cMiAaPQ $dMiAaPQ
050     #4$aBF575.P9 .N58 2019
082     0#$a303.3/85
100     1#$aNittle, Nadra.
245     10$aRecognizing Microaggressions.
250     ##$a1st ed.
264     #1$aNew York, NY : $bEnslow Publishing, LLC, $c2018.
264     #4$c©2019.
300     ##$a1 online resource (82 pages)
336     ##$atext $btxt $2rdacontent
337     ##$acomputer $bc $2rdamedia
338     ##$aonline resource $bcr $2rdacarrier
490     1#$aRacial Literacy Ser.
505     0#$aCover -- Title Page -- Copyright -- Contents -- Introduction -- CHAPTER 1 --
588     ##$aDescription based on publisher supplied metadata and other sources.
650     #0$aMicroaggressions.
650     #0$aPrejudices.
650     #0$aDiscrimination.
655     #4$aElectronic books.
776     08$iPrint version: $aNittle, Nadra $tRecognizing Microaggressions $dNew York, NY
830     #0$aRacial Literacy Ser.
```

# Data: Selected MARC Fields

| Label | Field | Subfield |
|-------|-------|----------|
| bibid | 035 | a |
| Title | 245 | a |
| Publication Date | 260 | c |
| Description | 520 | a |
| Topical Subject Terms | 650 | a |

# Methods: Data Acquisition

1. **Obtained** permission **from Dean of Libraries and Ex Libris company**

2. Exported 30 .tr.gz files from Alma ILS to Box: 11 hours

3. **Added files to SMU's** high performance **computer**
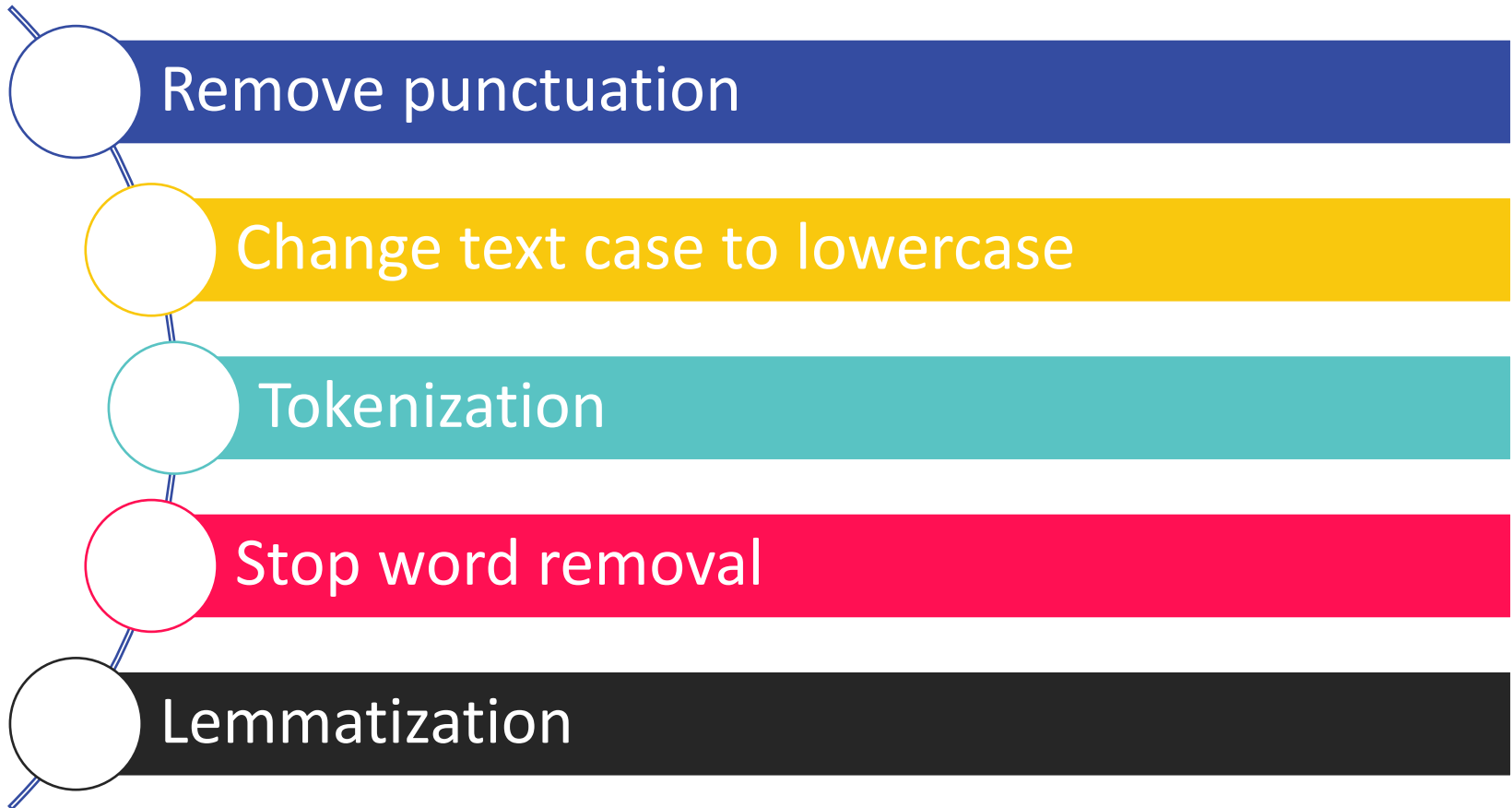
4. Decompressed into 19.8GB of XML files

5. **XML Parsing** with ElemenTree **XML API**

# Methods: Preprocessing

Remove punctuation

Change text case to lowercase

Tokenization

Stop word removal

Lemmatization

# Methods: EDA

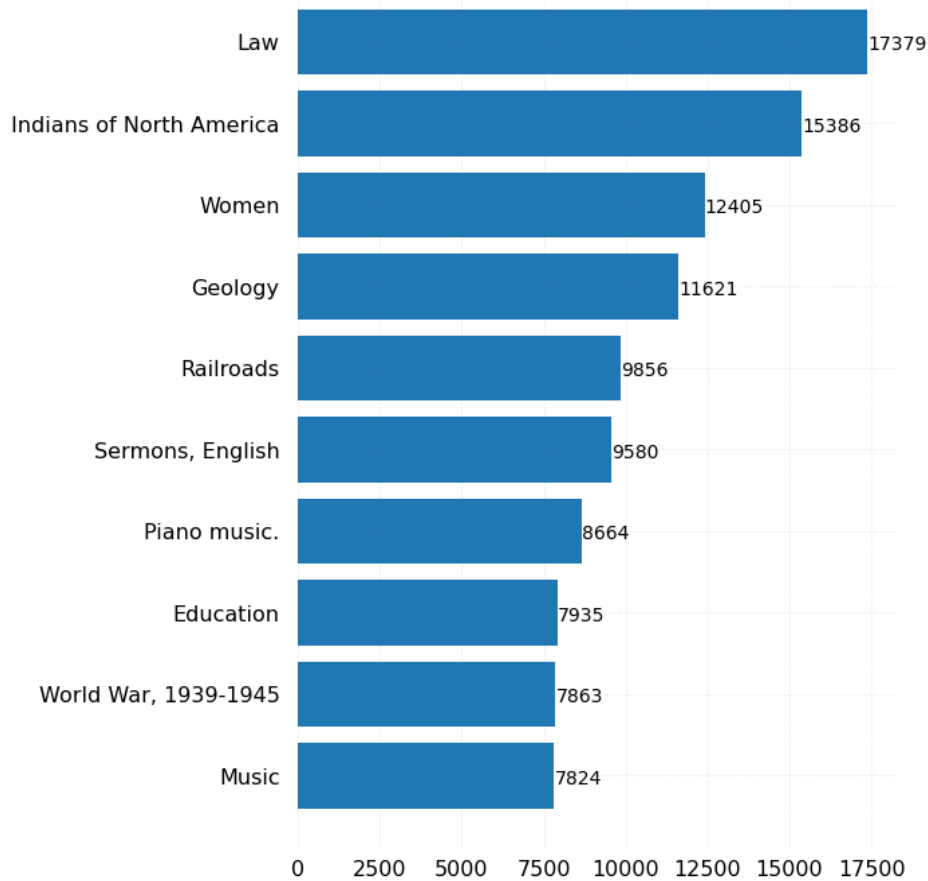| Field | Available Values | Percent of Missing |
|---|---|---|
| bibid | 5,475,642 | 0.11% |
| Title | 5,475,631 | 0.11% |
| Publication Date | 2,291,871 | 58.19% |
| Description | 721,454 | 86.84% |
| Topical Term Subject | 3,022,612 | 44.86% |

# Methods: EDA



Histogram of Frequency of Publication Year

# Methods: EDA



Top 10 Topical Subject Headings in the First Position of MARC Record

# Methods: Detoxify Model

- Developed by Laura Hanu, Unitary AI

- Primary model used to identify toxic language in the SMU catalog.

- Trained on Jigsaw Toxic Comment Classification Challenge Dataset.

# Results: Top Ten Titles

| bibid | Original Title | Toxicity Score |
|---|---|---|
| (OCoLC)1110021888 | Witches, Sl*ts, Feminists : | 0.9952282 |
| (CKB)4100000008484946 | ¬øEs tu jefe un gilipollas?. | 0.9951066 |
| (OCoLC)ocm34699730 | La parra, la perra y la porra / | 0.9943474 |
| (OCoLC)760279912 | Witches, wife beaters, and w***** : | 0.9937238 |
| (OCoLC)ocn708243813 | Witches, wife beaters, and w*****: | 0.9937238 |
| (YBPDDA)ebc869075 | Why are f****ts so afraid of f****ts? : | 0.99368894 |
| (OCoLC)ocn232648188 | Duck, you sucker | 0.9931779 |
| (OCoLC)ocm30835997 | Wild Dick, the Indian slayer ... | 0.992767 |
| (OCoLC)ocn953830144 | Stupid f***ing bird : | 0.99224806 |
| (MiAaPQ)EBC869075 | Why are f****ts so afraid of f****ts? | 0.992063 |

# Results: Top Five Descriptions

| bibid | Original Description | Toxicity Score |
|-------|---------------------|----------------|
| (CKB)371000000042 0489 | Do you ever find yourself thinking, how could you be so s*****, you look f**, or you're a h******* mother? Are you afraid people will find out you've fooled them into thinking you're competent? If you're guilty of expressing these types of discouraging messages, then you have a b**** in your head. This self-critical behavior can wreak havoc with your life-it can keep you from getting the love you want, the raise you deserve, or even a good night's sleep.Dr. Plumez began to notice a pattern with her patients being too hard on themselves. She found that gentler approaches didn't work, but when s... | 0.98363435 |
| (CKB)100000000038 0783 | Praise for Hoovers bestseller How to Work for an Idiot: ,Anyone who has to work should read How to Work for an Idiot., USA Today ,Dr. Hoover recommends admitting that you are powerless over the jerks in your life. Otherwise, harboring all that resentment is like drinking a cup of poison and waiting for the jerk to d**,... | 0.9819403 |
| (YBPDDA)ebc11771 99 | Uncrossable rivers! Hospitable nomads! Rabid dogs! Marijuana fields! Hailstone flashfloods! Maidens on horseback! Underpants wrestling! Toxic mountain-top lakes! S***** westerners! And the mountain-biking - so much biking you're a*** will hurt just reading it. This is what happens when two young i***** set out on a gonzo ride across the wilds of Mongolia. | 0.98072183 |
| (OCoLC)ocm371150 33 | Sleuth Hap Collins and his black sidekick go after a gang of gay-bashers. The gang makes videotapes as they r***, t****** and m****** homosexuals. The setting is Texas. | 0.9784098 |
| (TxDaM)3791072-smudb | "A g***b*** directed by the girl getting b*****. A h**** wrestler breaking into a teammate's locker to j*** o**. Two f*** in love drinking 40's and f****** on a rooftop. A burlesque start getting seduced by her real life girlfriend. A p*** star pour milk all over her h****. Yes, the queer world is as hard to define as it is to ignore--and these foxes are here to give you a VIP pass to the whole scene. | 0.9650895 |

DataScience@SMU

# Results

- 1 in 3000 titles and descriptions scored above a .75.

- 73 out of the top 100 English titles were truly harmful.*

# **Discussion**

- Identified most harmful and non-inclusive catalog records.

- Biggest challenge was figuring out how to work with an unlabeled set of data.

# Discussion

- Nuances in defining what language is harmful and non-inclusive.

- Not only an SMU problem. Other academic libraries can perform a similar evaluation on their catalog.

# Ethics

- Evolution of Language
- Trauma in Remediation Process
- Censorship
- Future Research Needs

# Conclusions

- **Strategic Goal:** Creating equitable and inclusive spaces

- **Finding the Needle in the Haystack**:
  - Harmful & non-inclusive text exists in catalogs
  - Machine learning tools can help

- **Opportunities**

# Questions

# Citations, page 1

1.Limbong, A., (2020) "Microaggressions are a big deal: How to talk them out and when to walk away." NPRLife Kit.https://www.npr.org/2020/06/08/872371063/microaggressions-are-a-big-deal-how-to-talk-them-out-and-when-to-walk-away. Accessed July 30, 2022.

2. Nadal, K., Issa,M., Leon, J., Meterko, V., Wideman,M., Wong, Y. "SexualOrientation Microaggressions: "Death by a Thousand Cuts" for Lesbian, Gay, and Bisexual Youth." Journalof LGBT Youth,8(3),

3.Interview with former student, confidential identity. Date: around June 1, 2022, location: SMU campus.

4. Julie Minikel-Lacocque. (2013). Racism, College, and the Power of Words: Racial Microaggressions Reconsidered. *American Educational Research Journal*, *50*(3), 432–465. https://doi.org/10.3102/0002831212468048

5. Blithe, & Elliott, M. (2020). Gender inequality in the academy: microaggressions, work-life conflict, and academic rank. *Journal of Gender Studies*, *29*(7), 751–764. https://doi.org/10.1080/09589236.2019.1657004

6. Nadal, Griffin, K. E., Wong, Y., Hamit, S., & Rasmus, M. (2014). The Impact of Racial Microaggressions on Mental Health: Counseling Implications for Clients of Color. *Journal of Counseling and Development*, *92*(1), 57–66. https://doi.org/10.1002/j.1556-6676.2014.00130.x

.University of Michigan Libraries. (2022). Remediation of Harmful Language in Library Metadata.https://www.lib.umich.edu/about-us/policies/remediation-harmful-language-library-metadata. Accessed June 23, 2022.

.SMU Libraries. (2019). Strategic Plan,2019-2024. https://www.smu.edu/libraries/about/dean/strategic. Accessed May 15, 2022.

7.Sanabria. S.H. personal communication. Email.July 18, 2022.

8.George, K., Grant, E., Kellett,C., & Pettitt, K. (2021). A Path for Moving Forward with Local Changes to the Library of Congress Subject Heading "Illegal Aliens".Library Resources & TechnicalServices; Library Resources & Technical Services, 65(3), 84-95.https://doi.org/10.5860/lrts.65n3.84

9.American Libraries (2021). "Library of Congress Changes Illegal Aliens Subject Heading". https://americanlibrariesmagazine.org/blogs/the-scoop/library-of-congress-changes-illegal-aliens-subject-heading. Accessed July17, 2022.

10.Fox, V. (2018). " Creating Change in the Cataloging Lab: Peer to Peer Review" Library Journal.https://www.libraryjournal.com/story/creating-change-in-the-cataloging-lab-peer-to-peer-review.Accessed June 23, 2022.

11.Cataloging Lab.(2022).https://cataloginglab.org/.Accessed June 23, 2022.

12.Hardesty, J. L., & Nolan, A. (2021). Mitigating Bias in Metadata: A Use Case Using Homosaurus Linked Data.Information Technology and Libraries; Information Technology and Libraries, 40(3),1-14. https://doi.org/10.6017/ital.v40i3.13053

13.Billey, A., Drabinski, E., & Roberto, K. R. (2014). What's Gender Got to Do with It? A Critique of RDA 9.7.Cataloging & Classification Quarterly, 52(4), 412-421.https://doi.org/10.1080/01639374.2014.882465

14.Wright, K. (2019). Archival interventions and the language we use.Archival Science, 19(4), 331-348.https://doi.org/10.1007/s10502-019-09306-y

15.McFadden, S., & Venker Weidenbenner, J. (2010). Collaborative Tagging: Traditional Cataloging Meets the "Wisdomof Crowds".The Serials Librarian, 58(1-4), 55-60.https://doi.org/10.1080/03615261003623021

DataScience@SMU

# Citations, page 2

16. Zhu, Z., Liang, J., Li, D., Yu, H., & Liu, G. (2019). Hot Topic Detection Based on a Refined TF-IDF Algorithm. IEEE Access, 7, 26996-27007. 10.1109/ACCESS.2019.2893980

17. Kim, D., Seo, D., Cho, S., & Kang, P. (2019). Multi-co-training for document classification using various document representations: TF–IDF, LDA, and Doc2Vec. Information Sciences, 477, 15-29. 10.1016/j.ins.2018.10.006

18. Sychev, O. A., & Penskoy, N. A. (2019). Method of lemmatizer selections in multiplexing lemmatization. IOP Conference Series. Materials Science and Engineering; IOP Conf.Ser.: Mater.Sci.Eng, 483(1), 12091. 10.1088/1757-899X/483/1/012091

19. Alfonseca, E., Garrido, G., Delort, J., & Peñas, A. (2013). WHAD: Wikipedia historical attributes data: Historical structured data extraction and vandalism detection from the Wikipedia edit history. Language Resources & Evaluation, 47(4), 1163-1190. https://doi.org/10.1007/s10579-013-9232-5

20. Library of Congress. (2009). What is a MARC record and why is it important? https://www.loc.gov/marc/umb/um01to06.html. Accessed July 29, 2022.

21. Library of Congress Network Development and MARC Standards Office. (2022). MARC 21 Format for Bibliographic Data. https://www.loc.gov/marc/bibliographic/. Accessed July 29, 2022

22. Manning, C. D., Raghavan, P., Schütze, Hinrich. (2009). Stemming and Lemmatization. Introduction to Information Retrieval. https://nlp.stanford.edu/IR-book/html/htmledition/stemming-and-lemmatization-1.html. Accessed September 9, 2022.

23. Dale, E., & Chall, J. S. (1948). A Formula for Predicting Readability. Educational Research Bulletin, 27(1), 11–28. http://www.jstor.org/stable/1473169

24. Shivam Bansal, C. (2022). textstat 0.7.3. https://pypi.org/project/textstat/. Accessed October 6, 2022.

25. Unitary. (2022). Detoxify 0.5.0. https://pypi.org/project/detoxify/. Accessed July 22, 2022.

26. Kaggle. (2017). Toxic Comment Classification Challenge. https://www.kaggle.com/c/jigsaw-toxic-comment-classification-challenge

27. Kaggle. (2019). Jigsaw Unintended Bias in Toxicity Classification. https://www.kaggle.com/c/jigsaw-unintended-bias-in-toxicity-classification/data

28. Hanu, L. (2020). How well can we detoxify comments online?. https://medium.com/unitary/how-well-can-we-detoxify-comments-online-bfffe5f716d7. Accessed October 6, 2022.

29. Howard, S. A., & Knowlton, S. A. (2018). Browsing through Bias: The Library of Congress Classification and Subject Headings for African American Studies and LGBTQIA Studies. Library Trends, 67(1), 74-88. https://doi.org/10.1353/lib.2018.0026

30. SMU Libraries Critical Cataloging Team, meeting with research team, June 23, 2022.

31. Deepanshi. (2022). Text Preprocessing in NLP with Python Codes. https://www.analyticsvidhya.com/blog/2021/06/text-preprocessing-in-nlp-with-python-codes/. Accessed July 22, 2022.

32. Planned Parenthood. (2022). Transgender Identity Terms and Labels. https://www.plannedparenthood.org/learn/gender-identity/transgender/transgender-identity-terms-and-labels. Accessed October 6, 2022.

33. Steinmetz, K. (2014). Why it's best to avoid the word 'transgendered'. https://time.com/3630965/transgender-transgendered/. Accessed October 6, 2022.

DataScience@SMU