EMOJI BOOK RECOMMENDER

AGENDA

- Introduction
- Use Case
- NLP Task: Information Retrieval
- Emoji data preparation & keyword annotation
- Preparing the dataset of books
- Walkthrough of the pipeline
- Model Evaluation
- Demo
- Conclusion and Limitations
- Future Directions



INTRODUCTION

- Problem definition: How can we create a unique, interesting way to recommend books to the user?
- Details considered:
 - How can we make the program accessible and easy to use?
 - What should we consider when making recommendations?
 - Genre? → broader topic words
 - What are potential limitations of the approach?
 - Relying on book descriptions
 - How might we evaluate such a model?
 - Subjectivity of relevance scoring



USE CASE: WHY THIS PROJECT?

- In some areas, a user might not have fully articulated keywords in mind for a search
- One case: looking for something new to read!
 - Users might not know exactly what topic they're looking for
 - They might not know exactly what genre or age range they want
 - They may not have a specific title or author in mind
 - ...But they tend to have at least light preferences in some of these areas
- Our idea is to experiment with a fun, engaging method of input in a context where precision may not be the goal for the user

NLP TASK: Information Retrieval



- Open-minded use case
- Distinct from in-class discussion of searching for a restaurant



Other differences:

- User does not input an actual keyword — more interpretation on the part of the program
- Weighting and doubling emoji input

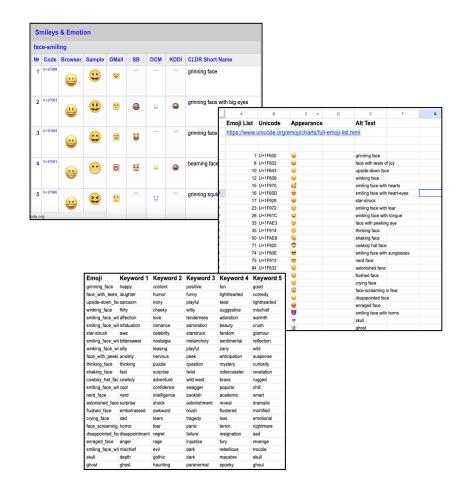
• • • • • • • • •



Datasets and Annotation

Emoji Dataset and Annotation

- Started with Unicode emoji list
- Pared down and added to spreadsheet, eliminated emoji that were:
 - Redundant (multiple smileys)
 - Not widely applicable in context (horoscope symbols)
- Added five keywords to describe each emoji
 - Started with list of Chat GPTgenerated keywords, then edited and replaced with our own tailored keywords
- Finally, formatted into TSV file—ready to work with!



EMOJI KEYWORDS



student	student	learning	school	education	study
teacher	teacher	education	school	instruction	classroom
judge	judge	law	justice	court	authority
farmer	farmer	agriculture	earth	gardening	nature

SHORT-TEXT LABEL: fed into back-end with Python emoji library, stored in keyword matrix **FIRST KEYWORD:** parallel to the short-text, the primary keyword and/or name of the emoji (after initial testing, we decided to weight this more heavily to get more topical books) **OTHER KEYWORDS:** related topics, settings, and descriptive words

WHERE DID WE GET THE BOOK INFORMATION? TRIAL AND ERROR

Try 1:

- Google Books
- Open Library
- Internet Archive

Queried using keywords derived from emoji inputs → 43,000 books

Stored data includes:

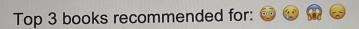
• Title, authors, publisher, description, etc.

Issues:

- Internet Archive included different media types not just books
- Open Library was missing book descriptions

20 (6) (7) (A)





Book 1: ('sad', 66.666666666666)

Book 3: ('Sad word 7', 25.0)

WHERE DID WE GET THE BOOK INFORMATION? TRIAL AND ERROR

Try 2: Kaggle!

 Kaggle data scraped from wonderbk.com (an online book store) which contains details of 100,000+ books

Stored data includes all the same information:

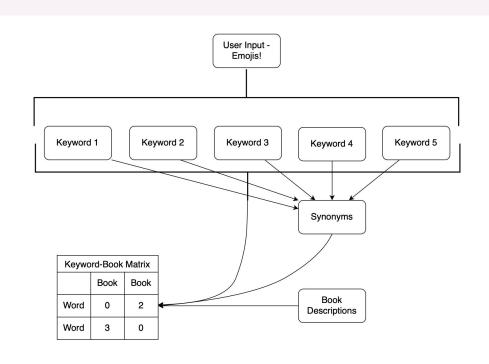
• Title, authors, publisher, description, etc.

Filtered out all books which didn't have a description field

Easy to parse and extract to make term-document matrix

BOOK DATASET: INCORPORATING INTO THE PIPELINE

- Book dataset file → book descriptions
- Into the pipeline! → Use book descriptions to count keywords and create a term-document matrix...



OUR PIPELINE



- Building the keyword-book matrix using Kaggle data
- User Interface
- Query Processing

BUILDING THE MATRIX

Book Title and Author

Twiliaht

War &

Dracula

		Farm	rwingrit	Peace	Diacula
	barnyard	2	0	0	0
	vampire	0	5	0	2
	Russia	1	0	2	1

Keyword

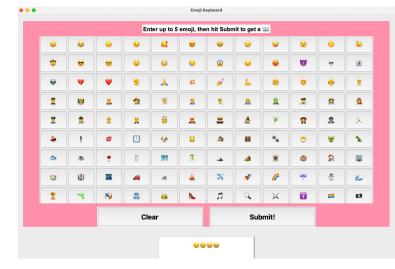
Large dataset led to complications

Animal

- Term-document matrix stores books per keyword (including synonyms)
 - Synonym support from NLTK WordNet
 - Normalized to adjust for description length (with raw counts, the search privileged books with longer descriptions)
 - Filter out books w/ empty descriptions

USER INTERFACE & MAIN LOOP

tkinter GUI:



- As simple as possible: clear instructions and visual design
- Interactive keyboard with display (Python emoji library)
- Simple buttons: option to reset input and submit
- No need to interact with terminal to operate
 - Once user submits their emoji selection, keyboard freezes and UI changes to a results page

PROCESSING BETWEEN USER QUERY AND MATRIX

- Emoji short-text associated to keywords
- Search through matrix file with user query
 - Matrix file is generated separately beforehand to cut down runtime
- Weight keywords
 - Book score = (count_book)*(count_query)
 - Prioritizes keywords with high frequency in both book and query
- Return sorted list of book scores

• • • • • • • • •

EVALUATION & CONCLUSIONS

MODEL EVALUATION

 Not always a single "correct" answer → hard to define accuracy

 Subjectivity: Relevance and accuracy depends on user interpretation

 Lack of labeled data for relevance → ideally we would get human annotators

SOFTWARE DEMO!







LIMITATIONS

- Difficulty in evaluating the model
 - True success is somewhat subjective (i.e. did they like the book?)
- Relying on descriptions of books that we didn't write ourselves → sometimes descriptions are not robust and/or topical
 - Need databases w/ descriptions, which may exclude certain kinds of books
- User can upload their own book database, but it's not completely "plug and play"

FUTURE DIRECTIONS



Generalizable to other media

Movies? Shows? Songs? Poems? Pieces of art?

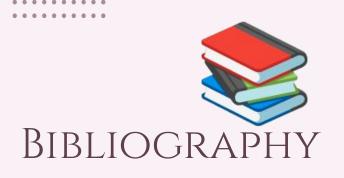
Filters for stronger user preferences

- Fiction vs. Nonfiction
- Length
- Excluding certain content

More accessibility features for the UI

- Used descriptive labels and buttons, but would ideally add more screen readers support
- Toggles for high-contrast emoji display
- Toggles to adjust the size of each emoji key for easier reading

Quality testing on a variety of devices



Dataset:

https://www.kaggle.com/datasets/elvi
nrustam/books-dataset?select=BooksD
atasetClean.csv

Recent Research:

https://www.nature.com/articles/s41598-O25-92286-O

- Presents a novel sentiment analysis approach that effectively incorporates emojis and emoticons in processing US airline tweets, using machine learning classifiers and BERT, achieving 92% accuracy—9% higher than state-of-the-art models. https://www.sciencedirect.com/science/article/pii/SO957417424O1O194
- Presents a Systematic Literature Review (SLR) of deep learning-based text summarization techniques (both extractive and abstractive) from 2014 to 2023, analyzing 73 studies to explore models, input representations, training strategies, evaluation metrics, datasets, and challenges in the field.
 https://www.researchgate.net/publication/334765586 Automated Genre Classification of Books Using Machine Learn ing and Natural Language Processing
- Proposes a machine learning-based genre classification method for books by transforming text into a feature matrix, reducing dimensionality using WordNet and Principal Component Analysis (PCA), and applying an AdaBoost classifier, achieving 92.88% accuracy by incorporating unlabeled data.

THANK YOU!

