



Internet Archive PDF

Huilin Chang, Yihnew Eshetu, Celeste Lemrow

Faculty Advisor : Professor Alvarado

Client : Internet Archive

Sponsor : Bryan Newbold

Presentation Outline

- Problem Statement
- Data Overview
- Data Pipeline
- Data Feature Engineering
- Exploratory Data Analysis
- Initial Models
- Conclusion

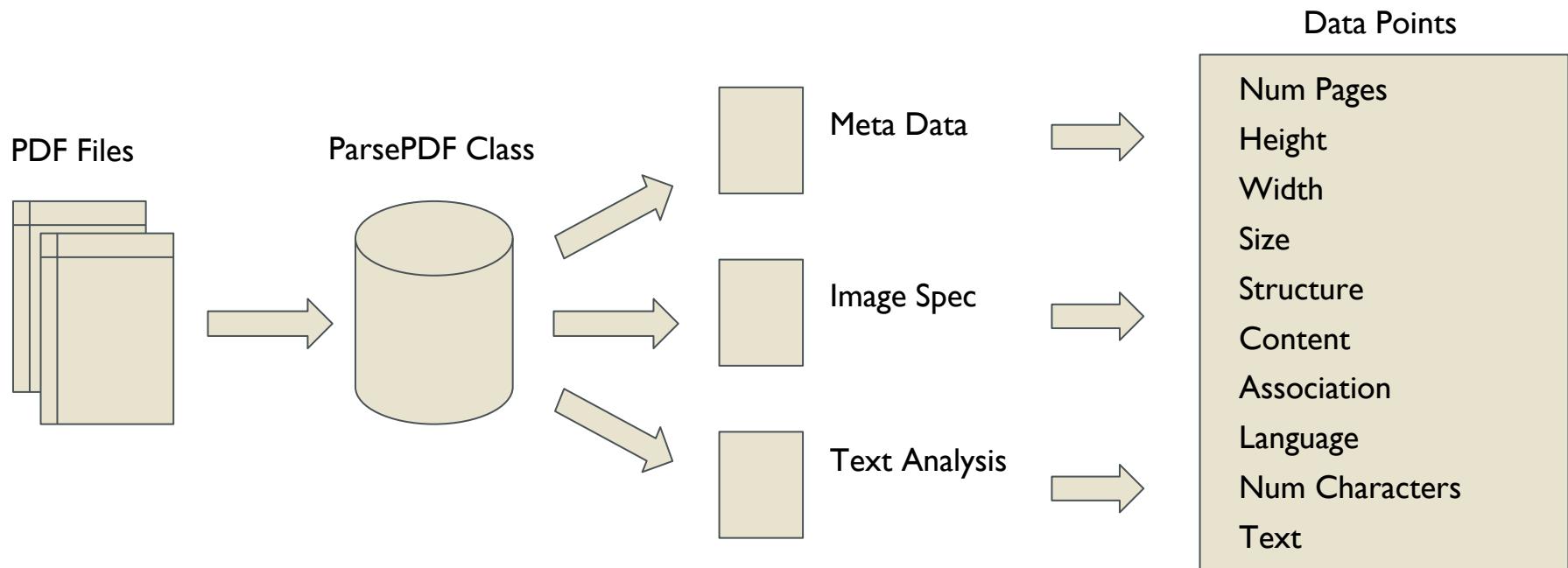
Problem Statement

- One of the Internet Archive's mission areas is “Universal Access to All Knowledge”, which includes collecting and providing access to the “scholarly web” -- research publications and datasets.
- Curation to accurately identify legitimate research publications is needed to help users find scholarly content
- An inclusive approach that accounts for diverse content, particularly from underrepresented geographic areas, groups, and content domains, is important to avoid excluding relevant content due to implicit bias and narrow criteria
- Our project aims to help this mission by implementing a fast PDF identification tool, which will score files on their likelihood of being a research publication

Data Overview

- 4 IA Training Datasets, consisting of 40k PDFs
 - Global Wayback Random - Random sampled PDFs from the Wayback Machine
 - Fatcat - A set of PDFs from the existing 'Fatcat' catalog of research papers
 - Fatcat Longtail Language - Papers from less-represented languages
 - Longtail - A set of PDFs created using heuristics (GROBID)
- Minor issues with data
 - Password Protected
 - Corrupted/Unparsable
- Plan to branch out further into IA content archives as well as other known sources of PDF scholarly documents

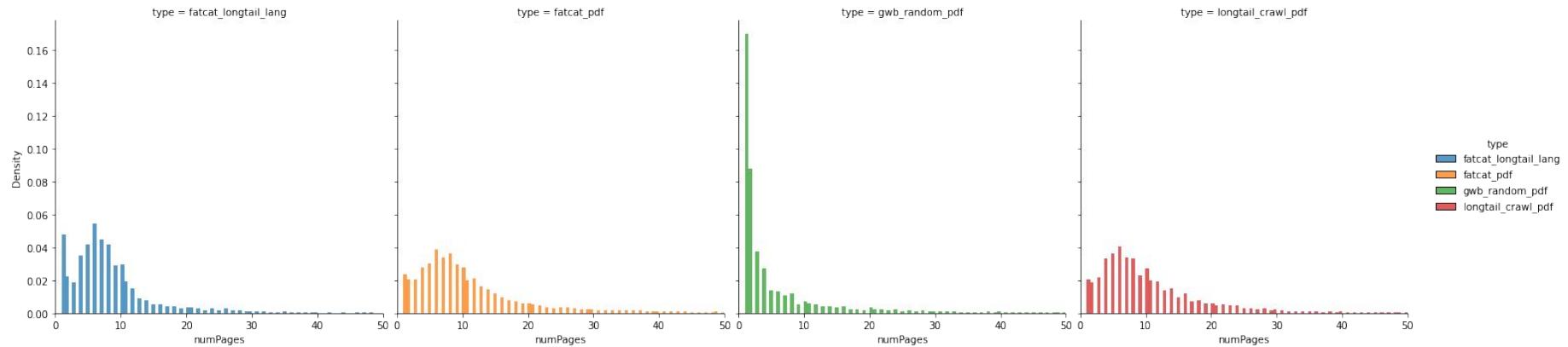
Data Pipeline



Data Feature Engineering

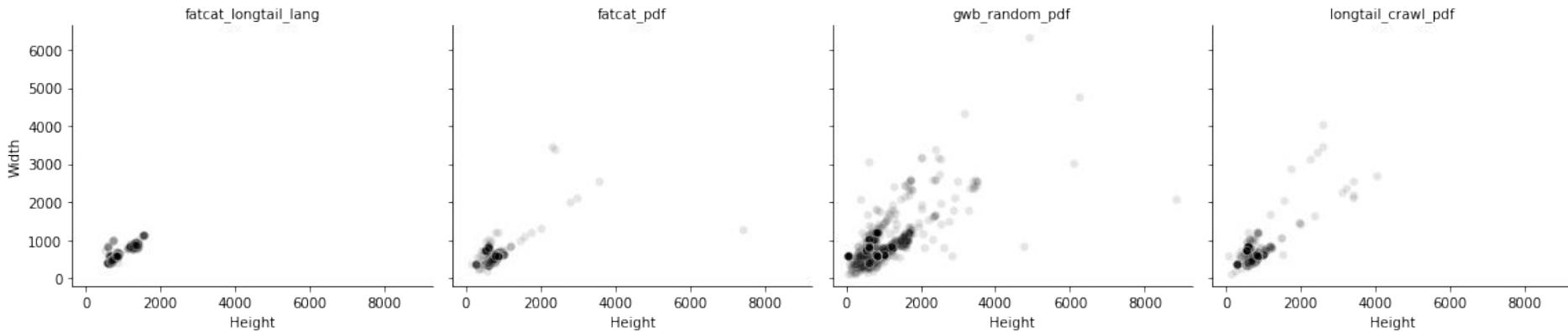
- **Structural Features:**
 - Number of Pages
 - Height
 - Width
 - Size
 - Number of Characters
- **Content Features:**
 - Language (English, Romance, Other)
 - Structure - words that signify structure of a research paper
 - Abstract, introduction, conclusion, reference, table of contents
 - Content - words that represent the content of a paper
 - Research, analyze, result, table, investigation, explain, theory, study, paper, data, perform
 - Association - words that represent association
 - Journal, association, organization, doi, university, school, board

Exploratory Data Analysis



Exploratory Data Analysis

Evaluation of document dimension by document type



Exploratory Data Analysis



Initial Models

- Bayes Model Averaging
- Topic Modeling (LDA)
- Decision Trees
- K-Means Clustering

Bayesian Model Averaging Results

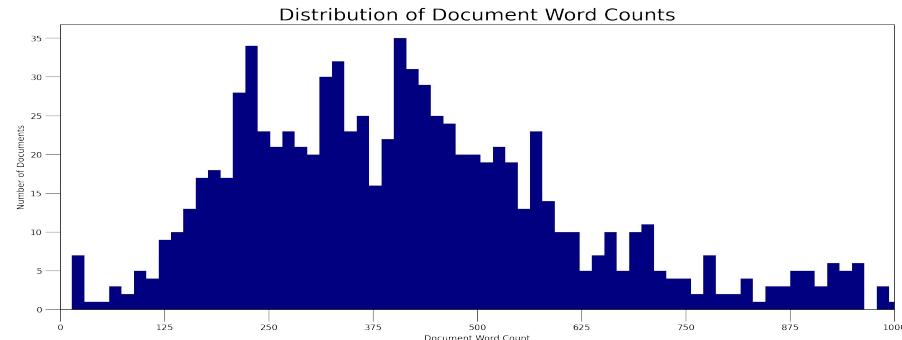
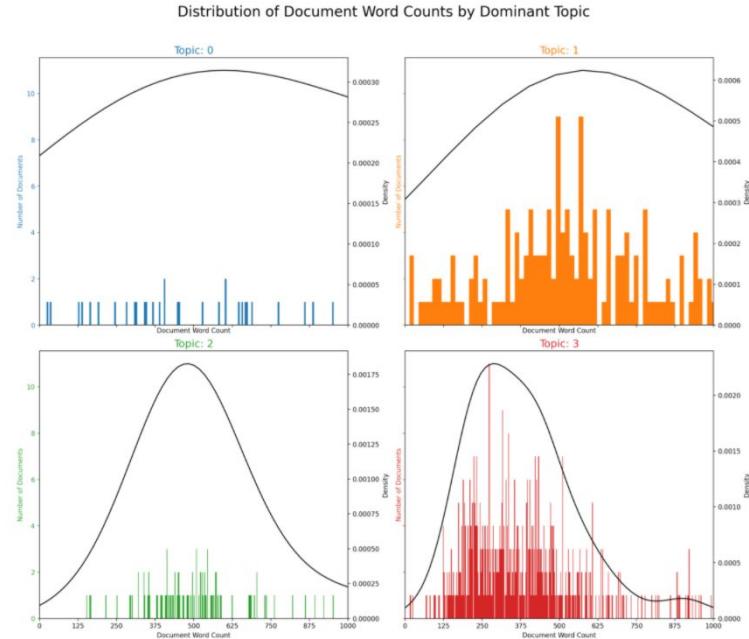
Model	Features	Likelihood
1	width, structure	6.20e-23
2	number of pages, width, structure, number of characters	1.03e-23
3	number of pages, structure, number of characters	5.71e-24

Model Prediction Comparison

Model	Features	Accuracy
1 - Top model from BMA	width, structure	75%
2 - 2nd best model from BMA	number of pages, width, structure, number of characters	78.3%
3 - Selected for comparison to assess content-focused features	structure, content, association	80%
4 - Full model	number of pages, height, width, dimension, structure, content, association, language, number of characters	77%
5 - Simple model (for baseline comparison)	structure	76.7%

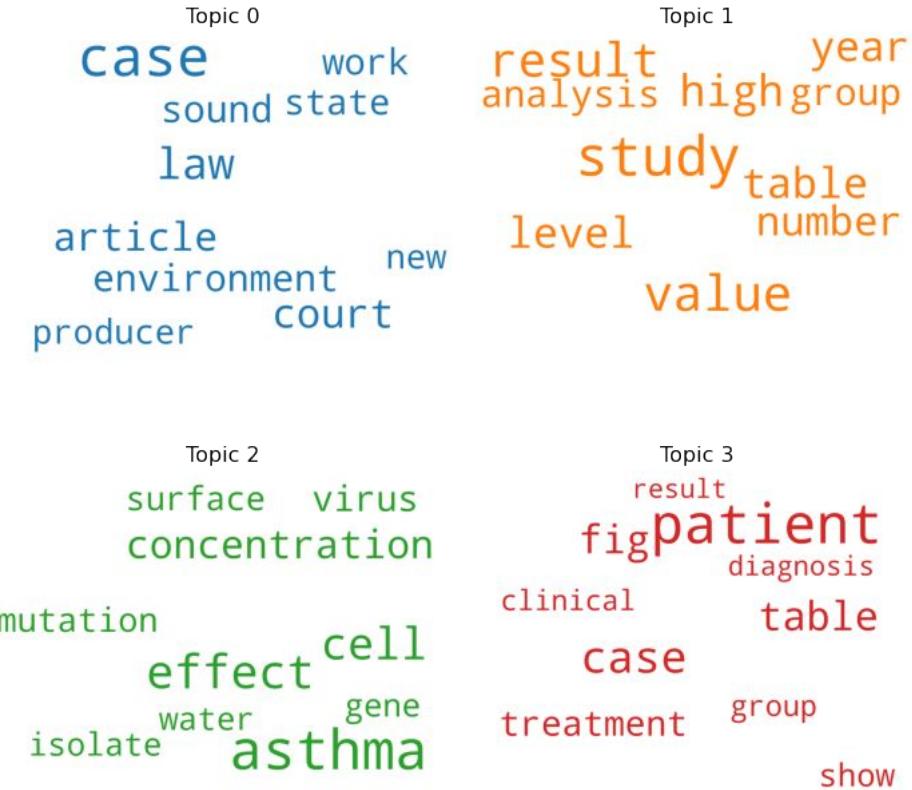
Topic Modeling

- Topic model based on LDA
 - Tokenize Sentences and Clean
 - Build the Bigram, Trigram Models and Lemmatize
 - Build the Topic Model
 - Dominant topic and percentage contribution in each document
 - Frequency distribution of word counts in documents

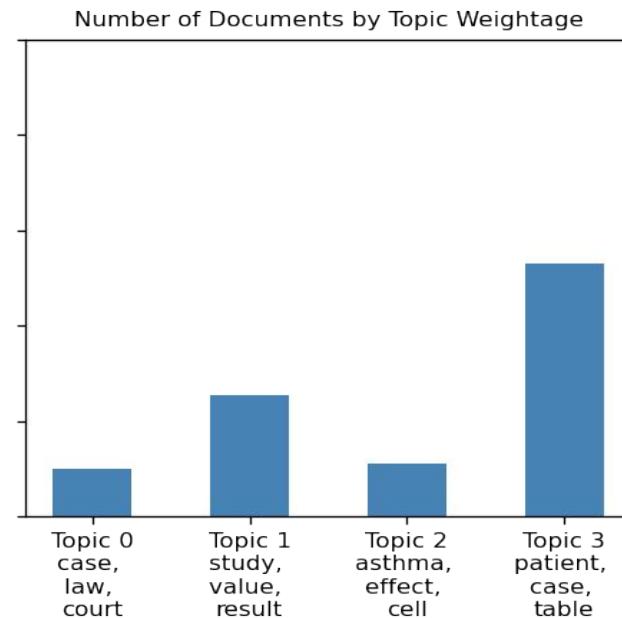
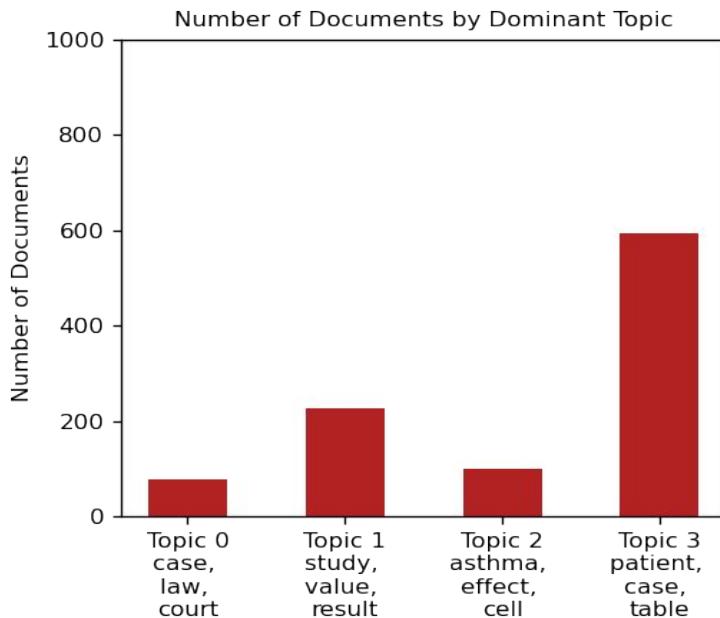


Topic Word Clouds

Top 10 Keywords of each Topic

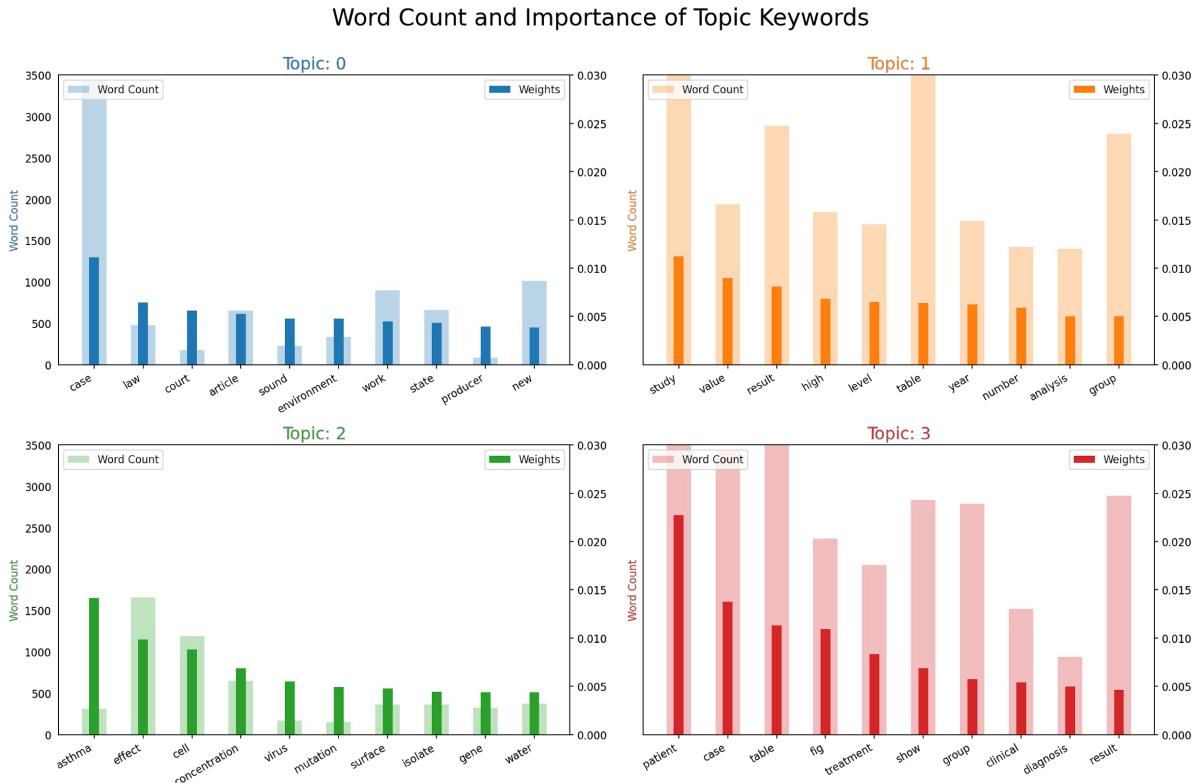


Top Topics



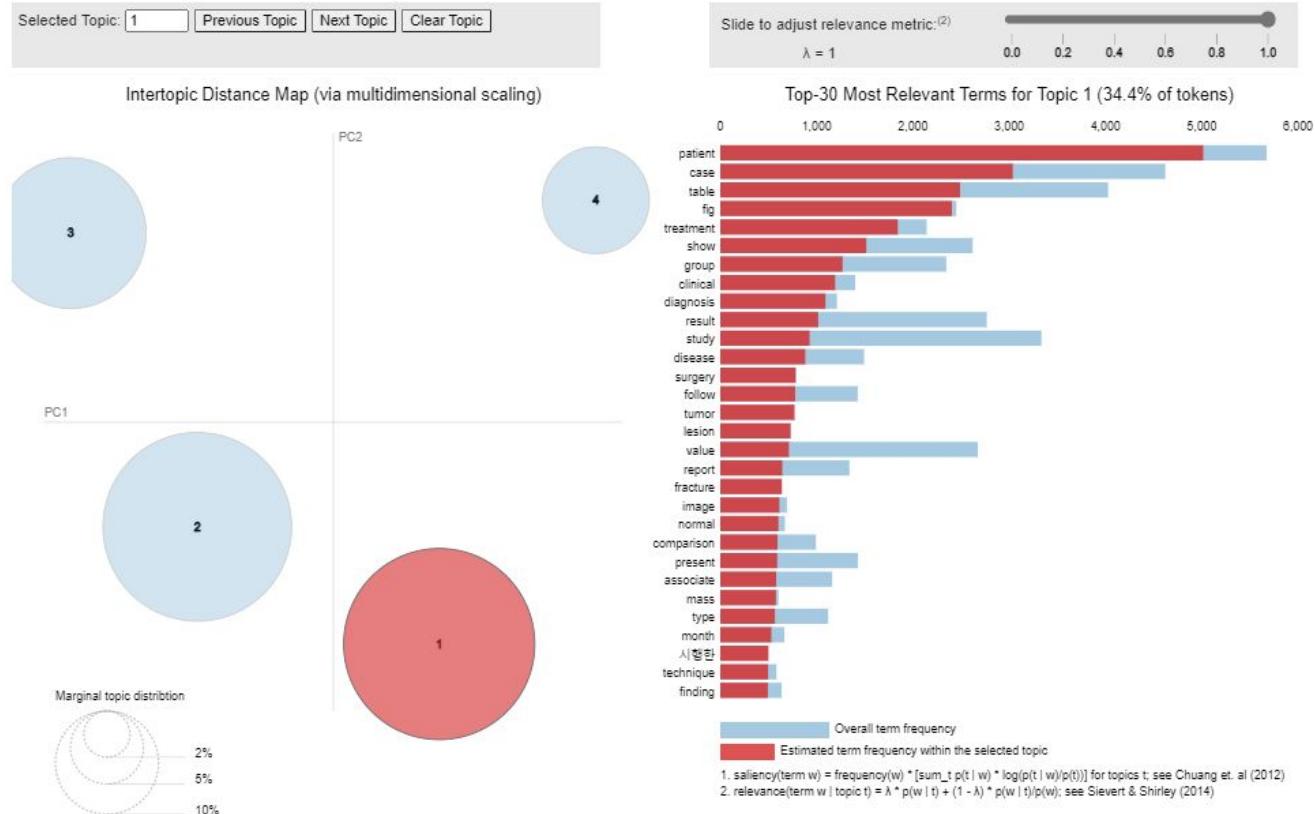
Word Count and Weighted Keywords

Frequency of Keywords in each Topic



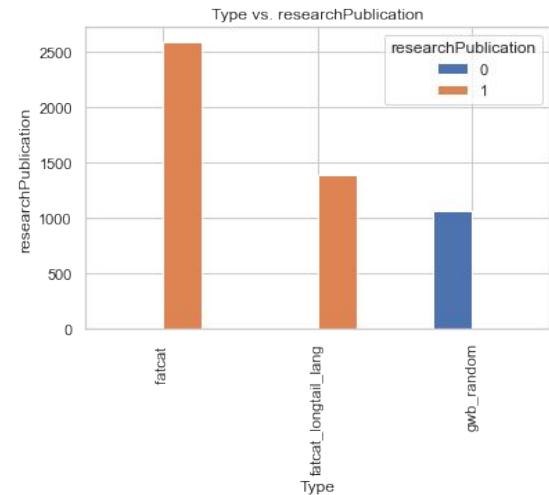
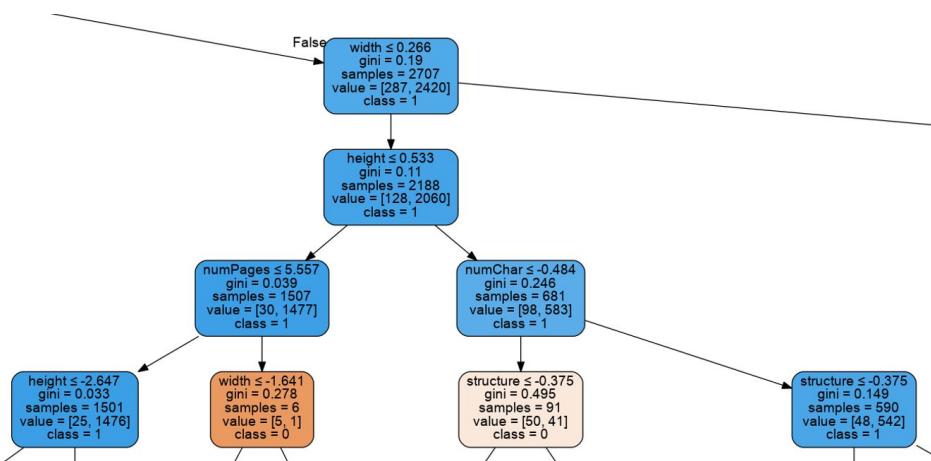
pyLDAvis

pyLDAVis Visualization of Topic Modeling



Decision Tree

- Nodes : Test for value of a certain attribute
- Gini index: how often a randomly chosen would be incorrectly identified.
- Blue : 1 Research Publication
- Orange : 0 Not Research Publication



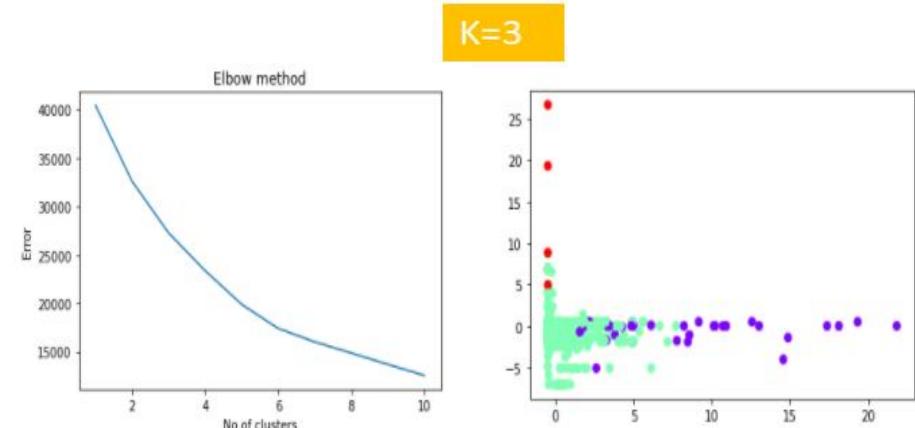
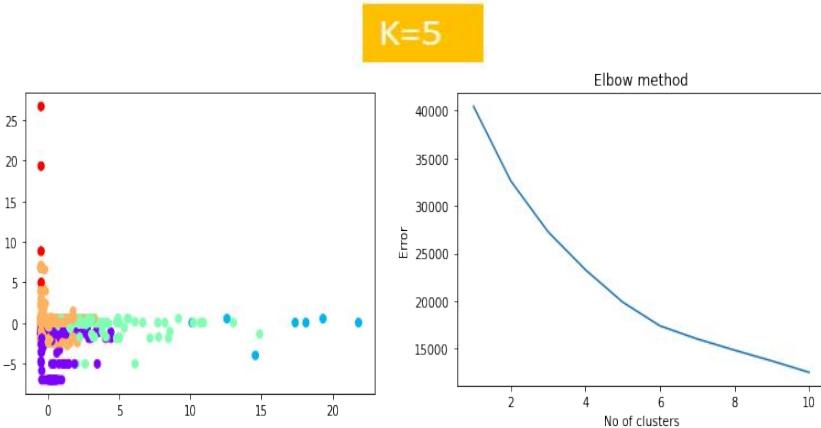
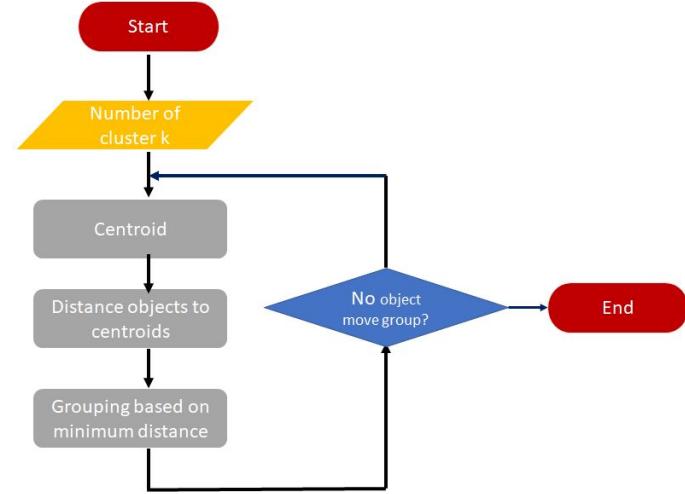
K-Means Clustering

Strengths

- Simple iterative method
- User provides “K”

Weaknesses

- Difficult to guess the correct “K”



Conclusions

- Document structural parameters may be powerful predictors of a research paper
- The “look” of a document can be leveraged for accurate prediction
- Language model content markers also can predict research content at a substantive level
- Multiple language pose a challenge -- require further consideration

Next Steps

- Expand structural and language model analysis to a larger set of documents and refine further
- Further integrate initial topic modeling results into the language model
- Additional parsing of individual content language model elements to determine comparative predictive power (for simpler, more efficient modeling without sacrificing accuracy)
- Experiment with language model analysis in different languages
- Build on document structure findings with a computer vision/ image model
- Develop methods to address “edge cases”
- Assess inclusivity to align with “long-tail” objective

Thank You!

Extra Slides

Sentence Topic

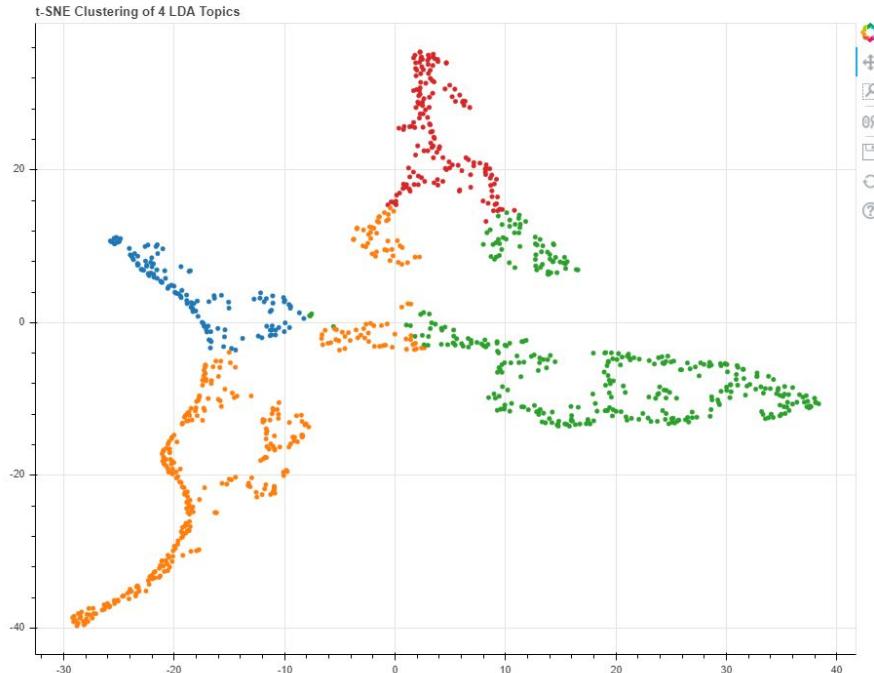
Words in documents are assigned color codes according to four topics

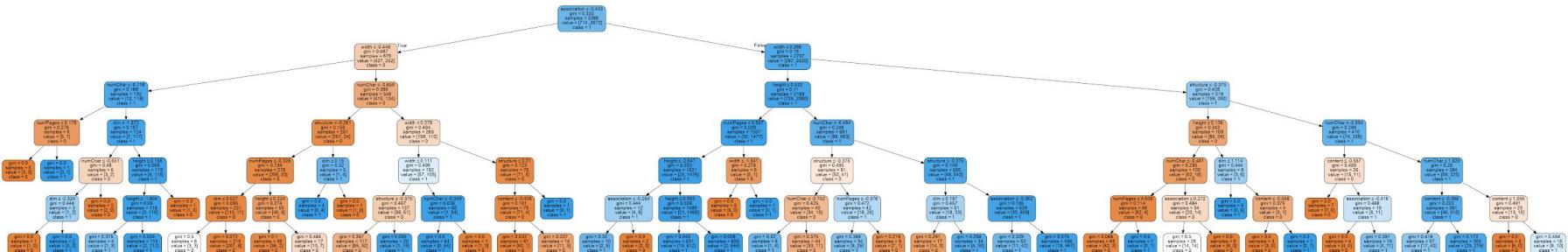
Sentence Topic Coloring for Documents: 0 to 11

Doc 0: according ageanalysis anterior_aspect aquarius aspect average axial calculate caucasian changeclao clinical cm ...
Doc 1: clinical conclusion distribution fig finding health korean leave medium_provide methodobtain open_access patient portion ...
Doc 2: ageanalysis change conclusion cooperation dependent describe entire female fig group handhealth human ...
Doc 3: analysis fig finding free healthy korean leave patient reference result role technique total ・・・
Doc 4: change comparison describe group normal old present reference related result structure study volume word ...
Doc 5: ageanalysis change clinical decrease female fig figure finding gender leave male normal obtain ...
Doc 6: clinical cmdistribution figure investigation medium_provide normal patient portion reference reproduction technique termapproach ...
Doc 7: ageclinical comparison conclusion correlation distribution female finding group health humaninvestigate medium_provide method ...
Doc 8: axial clinical microscopic normal present study case completely diagnosis protrude report ・・・ composeconnect ...
Doc 9: ageanalysis calculate change clinical conclusion female finding group health investigate mearmethodold ...
Doc 10: clinical comparison group humanpatient role structure study ・・・ con effect follow induce injection ...
Doc 11: change clinical fig leave malepatient reference study case diagnosis et_al experience managementreport ...

t-SNE Clustering Chart

t-SNE visualize the cluster of documents in a 2D space (t-distributed stochastic neighbor embedding)





Requirements for Project

Your presentation must have a minimum of 4 slides cover the material below.

- (a) Title with the names of your capstone team, faculty advisor, client, and sponsor.
- (b) Problem statement with a description of your problem and why it is important.
- (c) The sources and types of data you are using for your project.
- (d) Your data, **exploratory**, and **feature engineering** tasks to help you address your problem statement.

Data :

- 40k pdf files from 4 datasets

Exploratory :

- Graphs and Dashboard (**All of us can explore**)
- Bayes Model Averaging (model performance and parameters)

Feature Engineering :

- Meta data (csv)
- Text data (create csv containing file, text, and text language) **Yihnew**
- Images of first page of pdf
- Multiple languages (translate text to english)
- Top Modeling (**Huilin**)