# BAYESIAN CLASSIFICATION OF RESEARCH PAPERS

Huilin Chang(hc5hq) , Yihnew Eshetu (yte9pc), Celeste Lemrow (ctl7t)

DS6014 Project
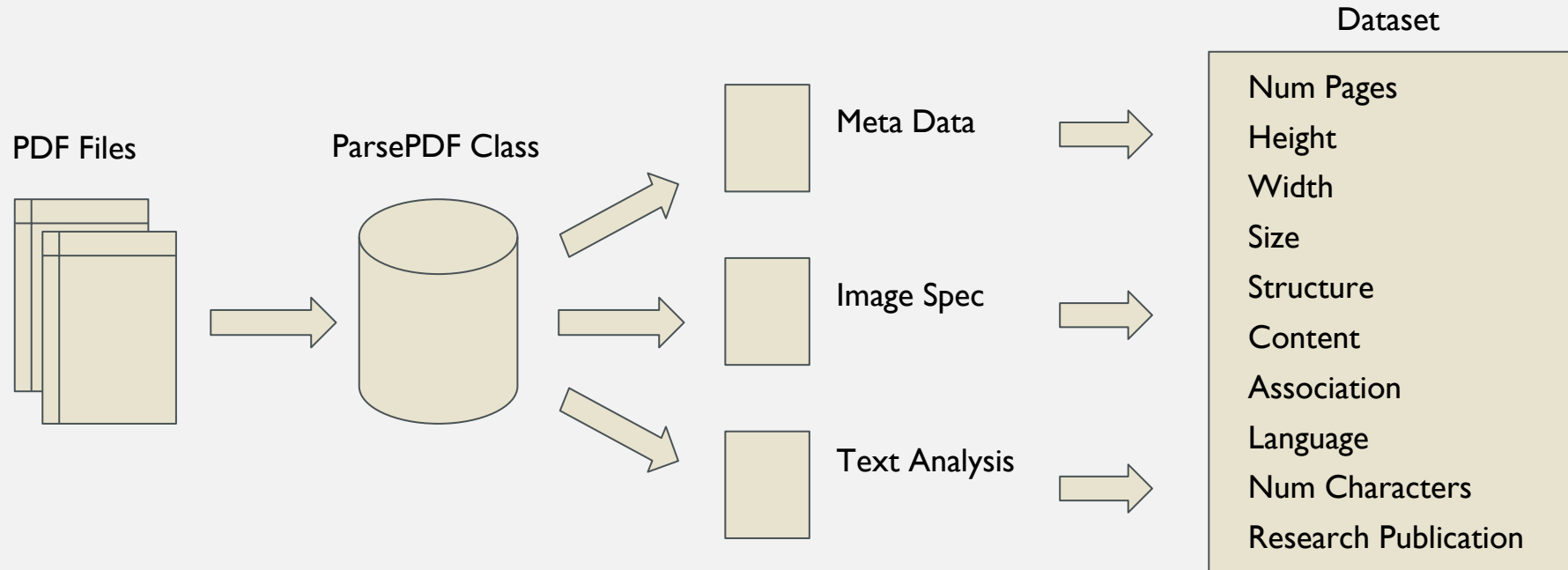
# TABLE OF CONTENTS

- Introduction

- Data Extraction

- Exploratory Analysis

- Methods

- Results

- Conclusions

# INTRODUCTION

- One of the Internet Archive's mission areas is "Universal Access to All Knowledge", which is an attempt to collect and provide access to the "scholarly web": the public record of research publications and datasets available on the world wide web.

- 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.

- Given the volume of PDF documents in the Internet Archive's (IA) repository, a classifier is needed to determine which are legitimate research documents

- Our group constructed several Bayesian logistic models with different feature combinations and compared results on accuracy of identification of research papers

# DATA EXTRACTION

PDF Files → ParsePDF Class → Meta Data / Image Spec / Text Analysis → Dataset

**Dataset**
- Num Pages
- Height
- Width
- Size
- Structure
- Content
- Association
- Language
- Num Characters
- Research Publication
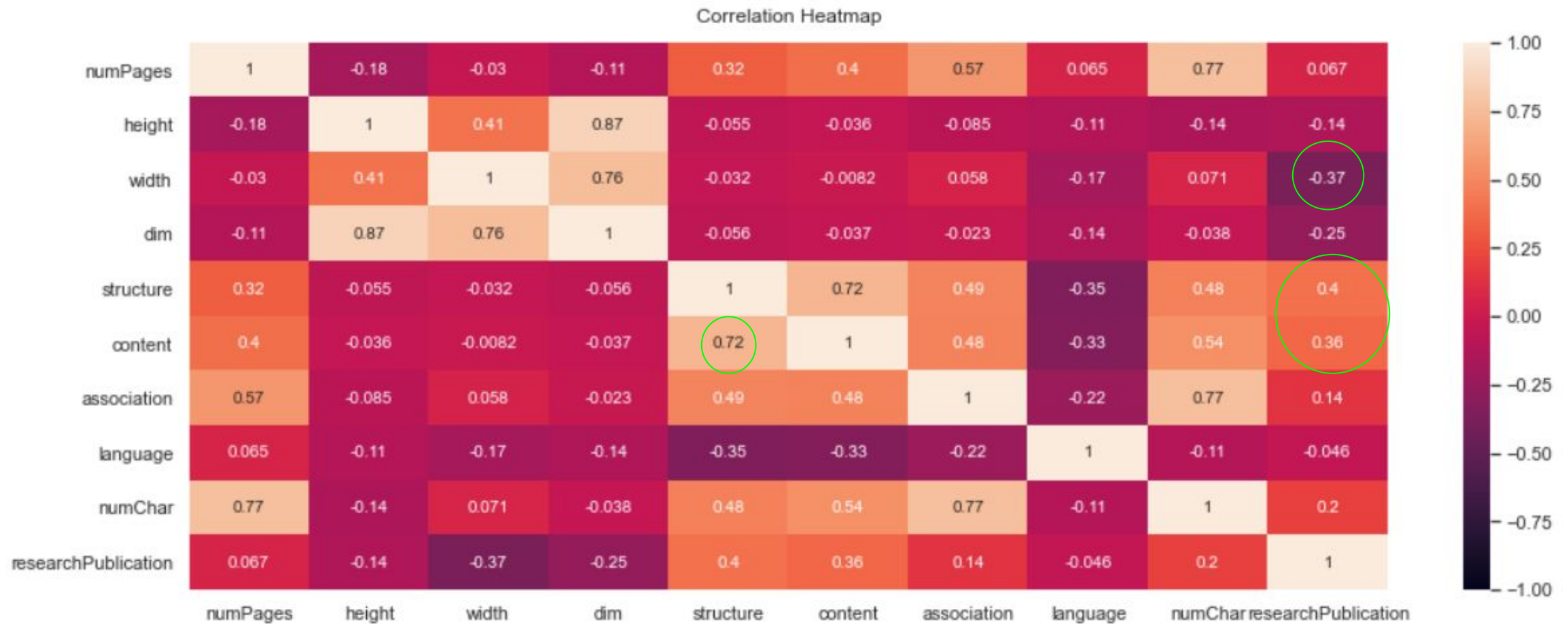
**Language:** English, Romance, and other
**Structure**: Words that represent the structure of a paper
        {abstract, introduction, conclusion, reference, table of content}
**Content**: Words that represent the content of a paper
        {research, analyze, result, table, investigation, explain, theory, study, paper, data, perform}
**Association**: Words that represent association
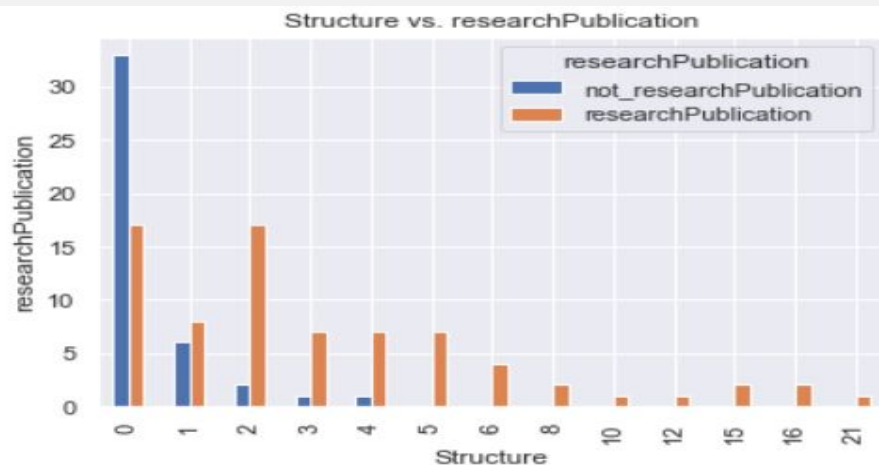        {journal, association, organization, doi, university, school, board}

# EXPLORATORY ANALYSIS

Correlation matrix : correlation between the research publication variable and selective features
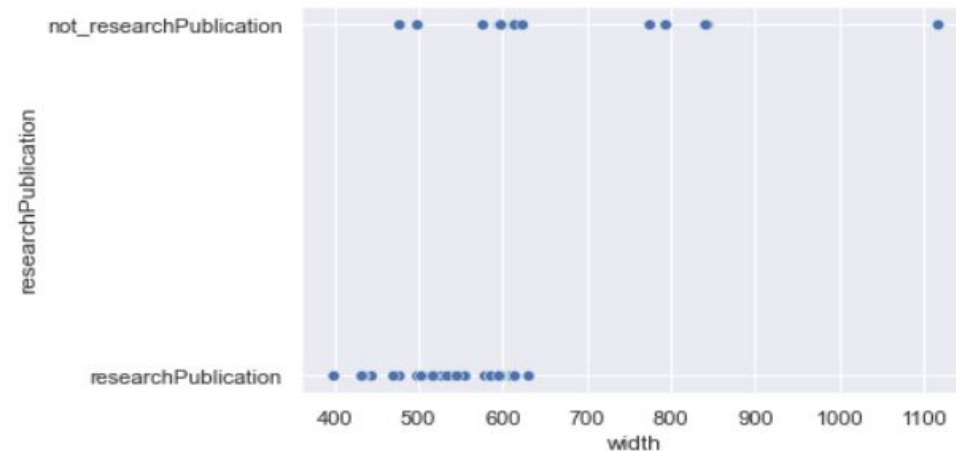


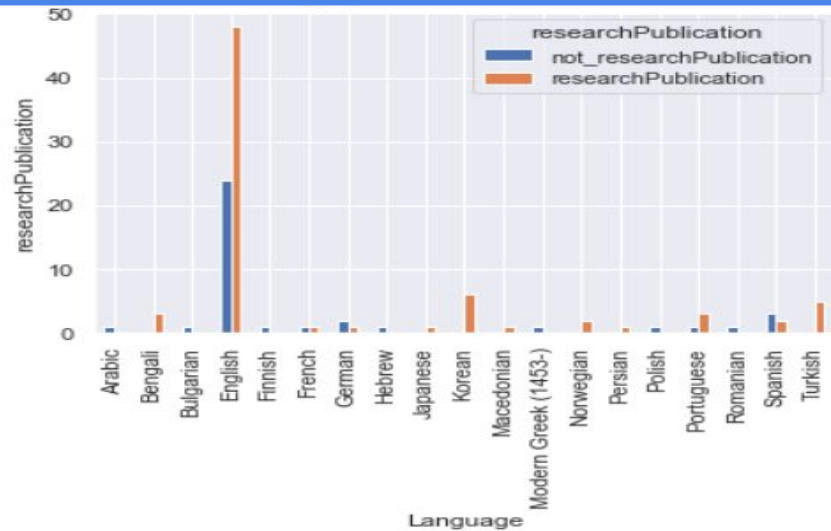Correlation Heatmap

# EXPLORATORY ANALYSIS


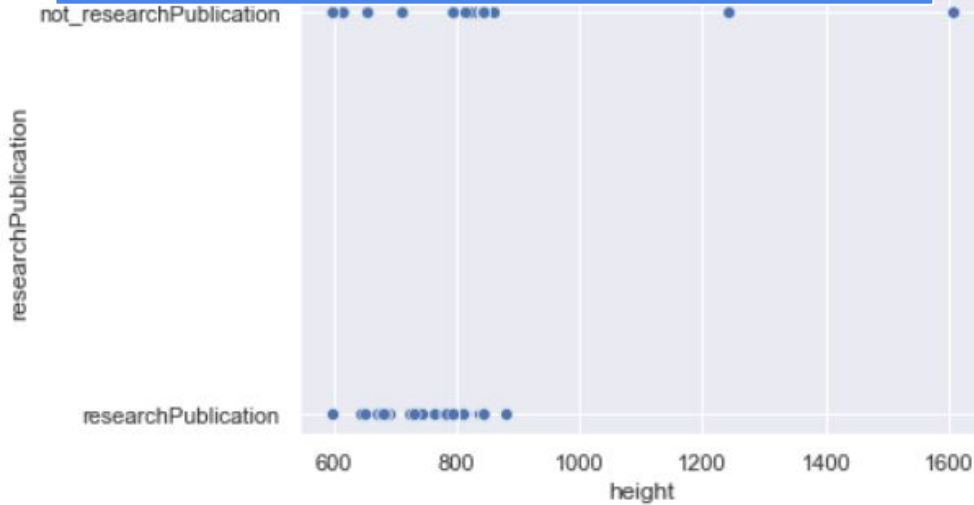
Structure vs. Research Publication



Width vs. Research Publication



Language vs. Research Publication



Height vs. Research Publication

# BAYESIAN STATISTICS

- Algorithm: Bayesian logistical regression, using PYMC3

$$p(\theta|D) = \frac{p(D|\theta)p(\theta)}{p(D)}$$

- Mathematical connection : the likelihood is the product of n Bernoulli trials,

$$\mathrm{II}_{i=1}^{n} p_i^{y} (1 - p_i) (1 - p_i)^{1-y_i} , \text{where } p_i = \frac{1}{1 + e^{-z_i}}$$
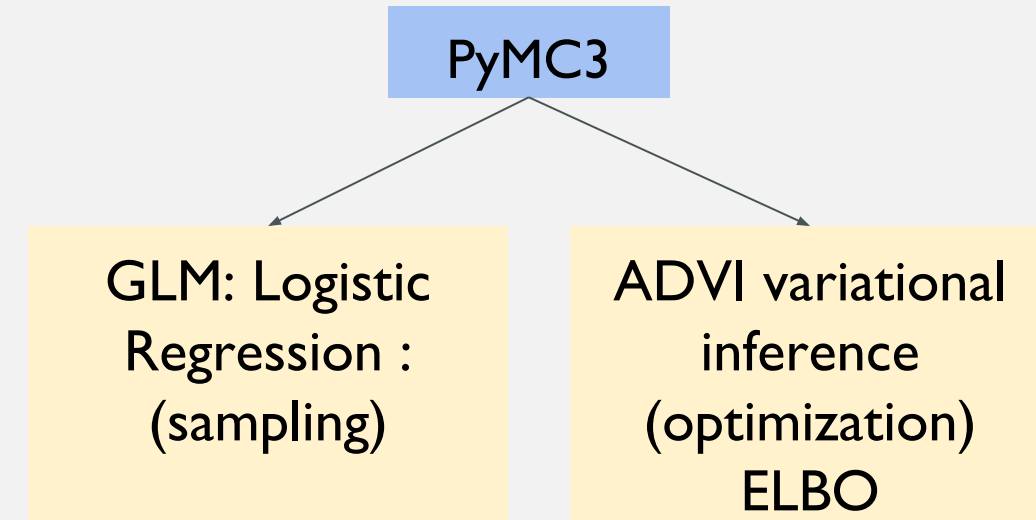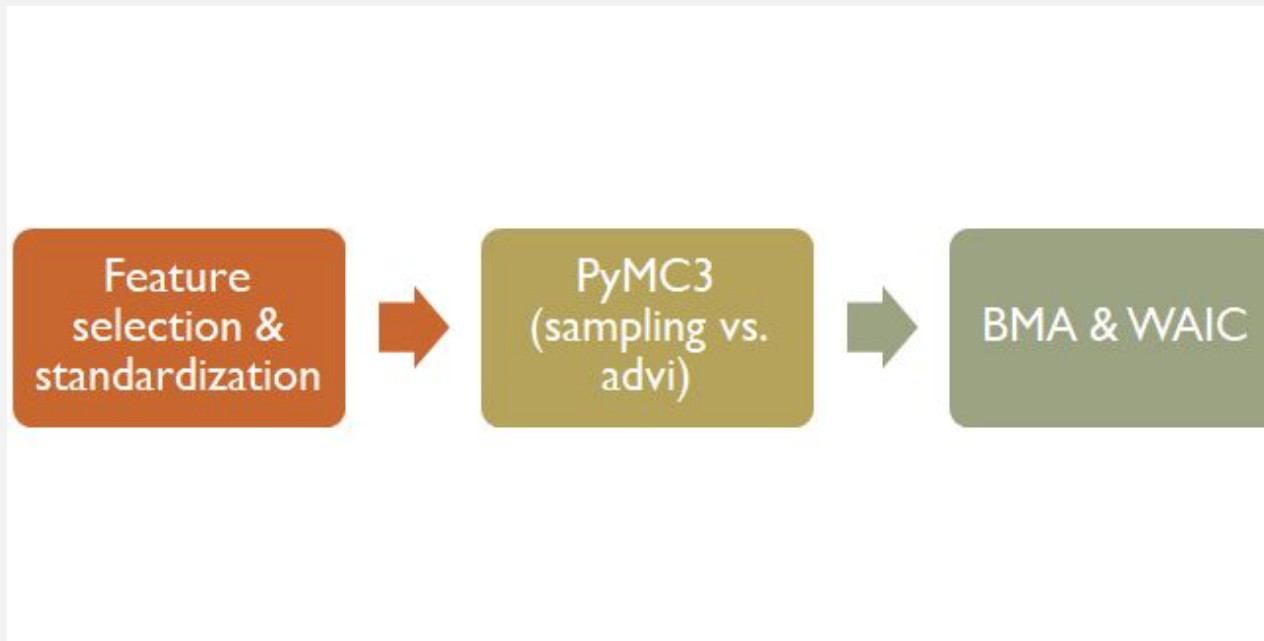
$$y_i = \beta_0 + \beta_1 (\text{numPages})_i + \beta_2 (\text{height})_{i+} \beta_3 (\text{width})_i + \beta_4 (\text{dim})_i + \beta_5 (\text{structure})_i + \beta_6 (\text{content})_i$$
$$+ \beta_7 (\text{association})_i + \beta_8 (\text{language})_i + \beta_9 (\text{numChar})_i$$

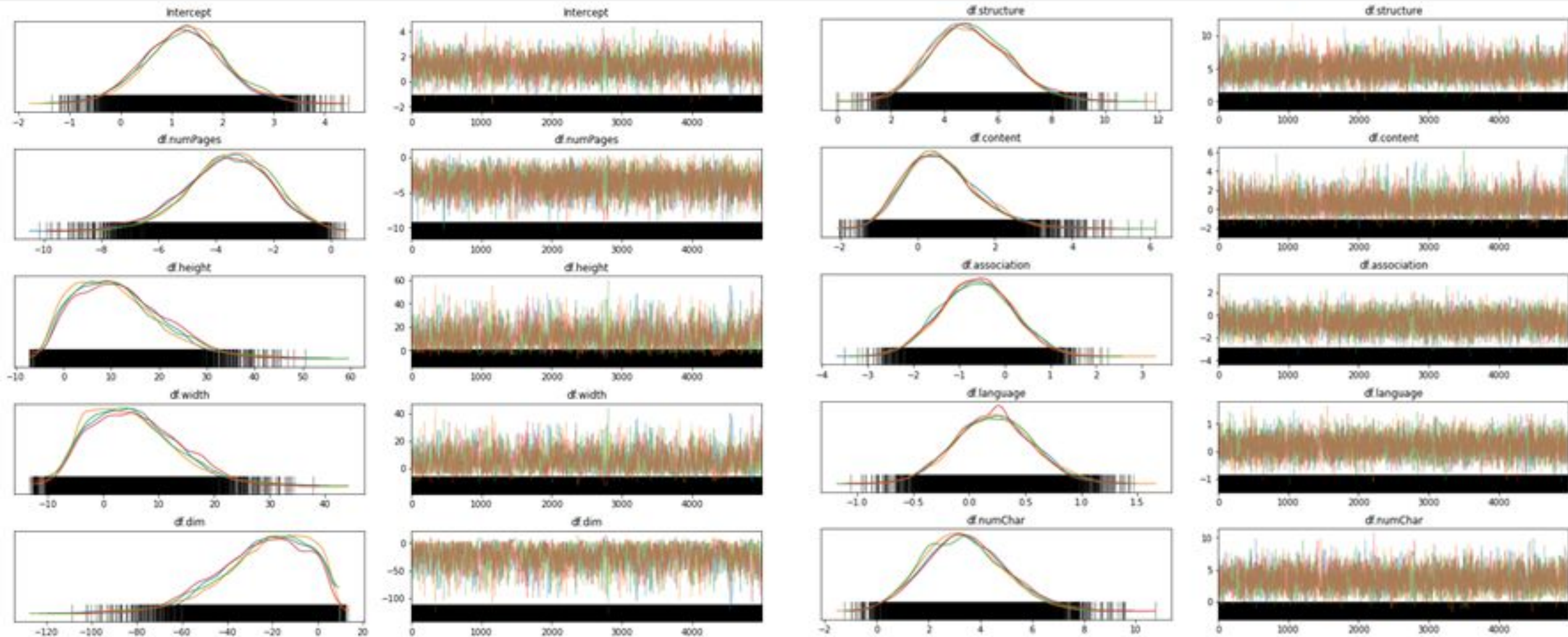Where $y_i$ = 1 if researchPublication and $y_i$ = 0 otherwise

- Priors : default $p(\theta) = N(0, 10^{12}I)$

# METHODS

- Total set of features considered: number of pages, height, width, dimensions of page, structure, content, association, language, number of characters

- How likely is it a **research publication** based on the selective features?

- Model comparison approach – compared different sets of features and accompanying accuracy. Given parameters for the capstone project, including speed, prioritizing a balance of the smallest number of features with acceptable accuracy is a key objective
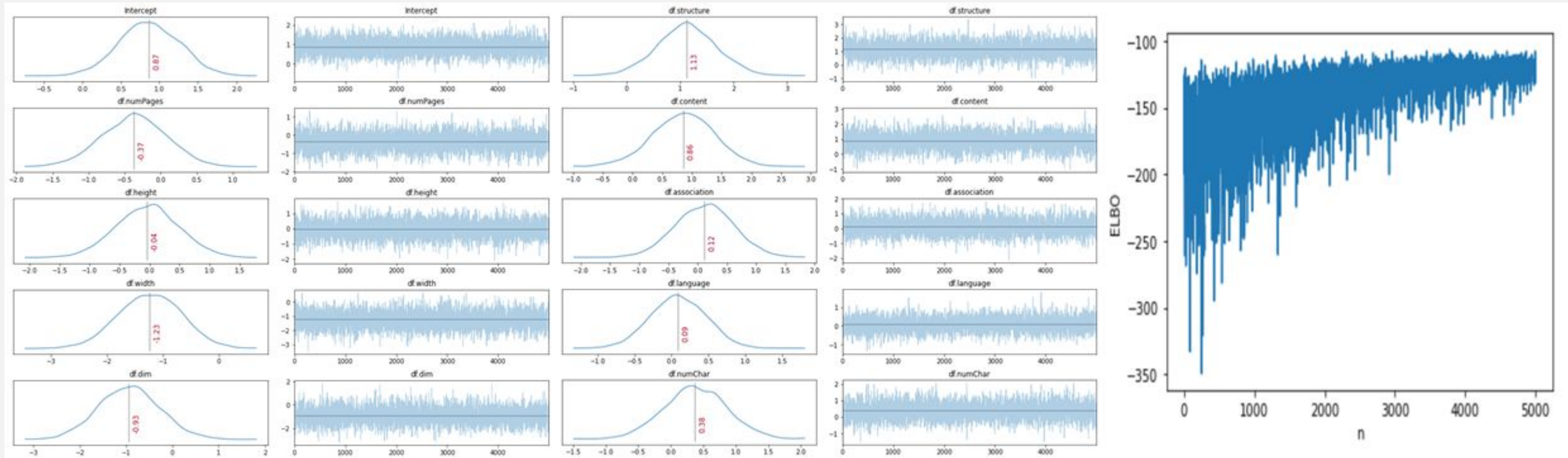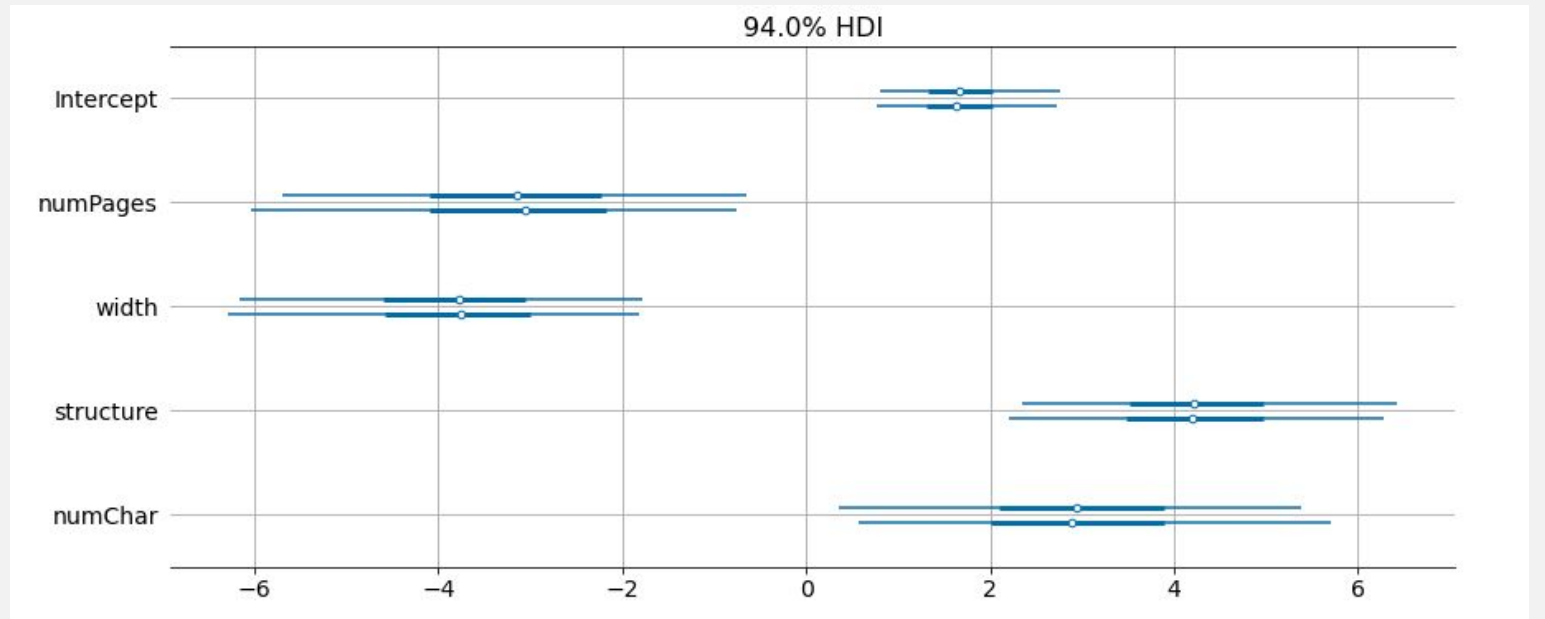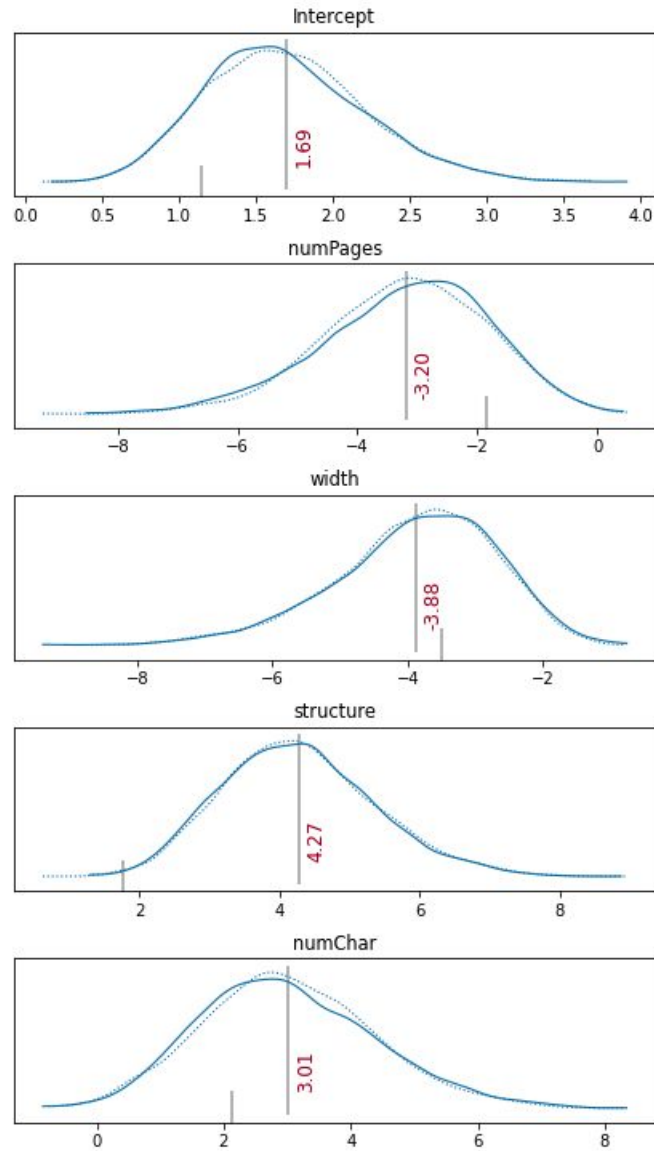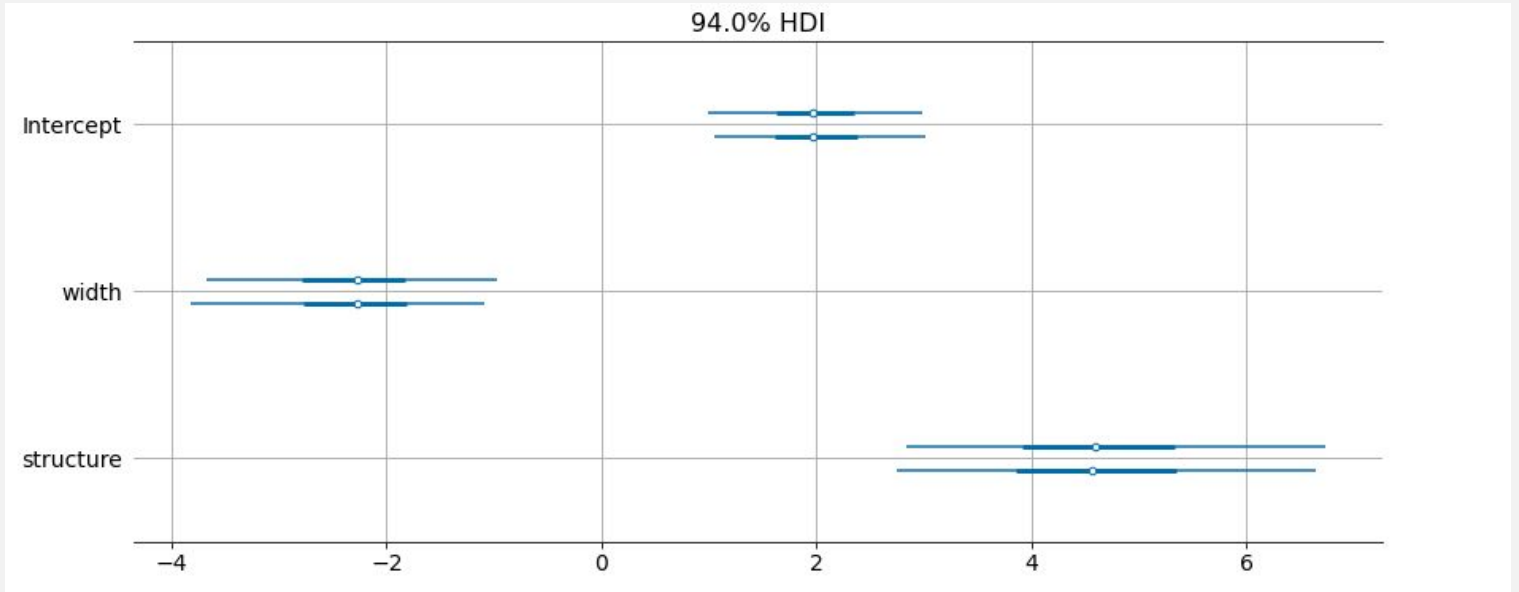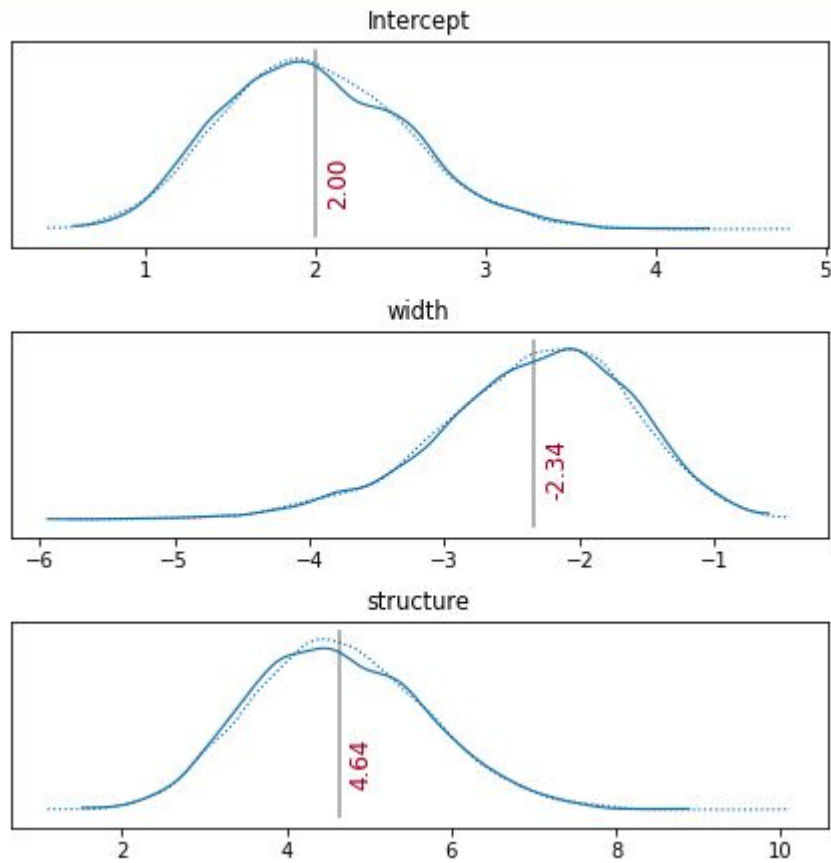
# RESULTS - Full Model (ADVI)

## Variational Inference: ADVI

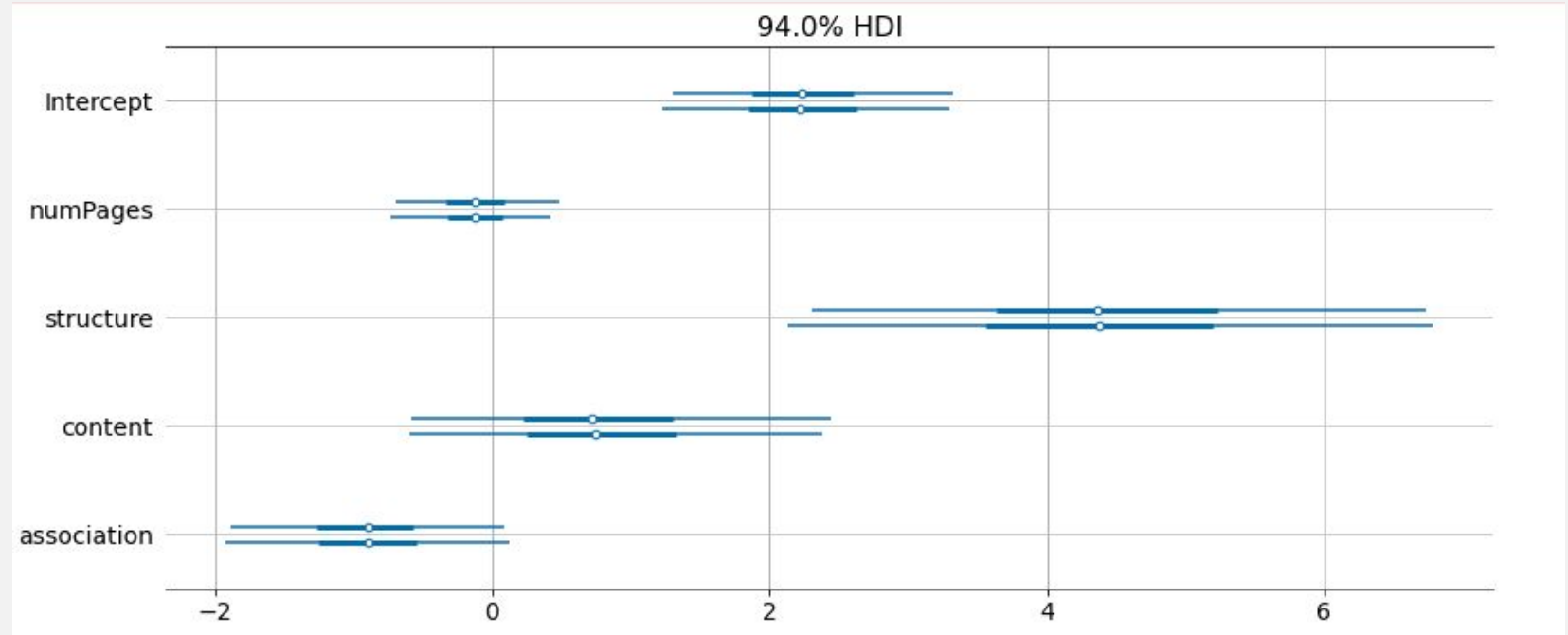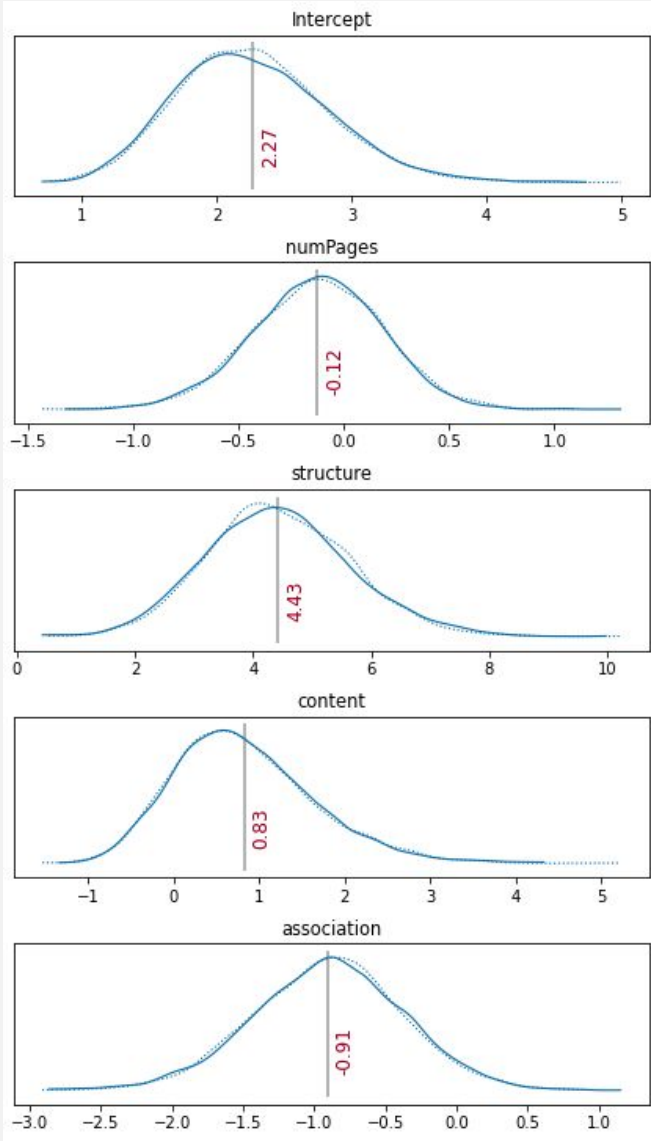# RESULTS - Model 1

# RESULTS - Model 2

# RESULTS - Model 3

# RESULTS - BMA - Top Models

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

# RESULTS - WAIC

| Model | Features | WAIC |
|---|---|---|
| 1 - Top model from BMA | number of pages, width, structure, number of characters | 95.0114 |
| 2 - 2nd best model from BMA | width, structure | 95.291 |
| 3 - Selected out of curiosity | structure, content, association | 119.888 |
| 4 - Full model | number of pages, height, width, dimension, structure, content, association, language, number of characters | 95.1717 |
| 5 - Simple model (also selected out of curiosity) | structure | 95.1717 |

# PREDICTION COMPARISON

| Model | Features | Accuracy |
|-------|----------|----------|
| 1 - Top model from BMA | number of pages, width, structure, number of characters | 78.3% |
| 2 - 2nd best model from BMA | width, structure | 75% |
| 3 - Selected out of curiosity | structure, content, association | 80% |
| 4 - Full model | number of pages, height, width, dimension, structure, content, association, language, number of characters | 77% |
| 5 - Simple model (also selected out of curiosity) | structure | 76.7% |

# CONCLUSIONS

- Key findings:
  - Top feature - structure
  - Additional key features - number of pages, width
  - Format matters, in addition to content
  - Simple model may be reasonable, given comparable accuracy to others, to prioritize speed

- Future work:
  - Cross-validation for further comparison
  - Additional investigation into models focused on the language model (structure, content, association)
  - Use LDA dimension reduction to assess whether it generates improved results
  - Expand analysis to larger dataset
  - Address non-English language factors in the model