# Classifying Text by Time Period

Zachary Kelly, Sasha Casada, & Nathan Le

# GOALS & MOTIVATIONS
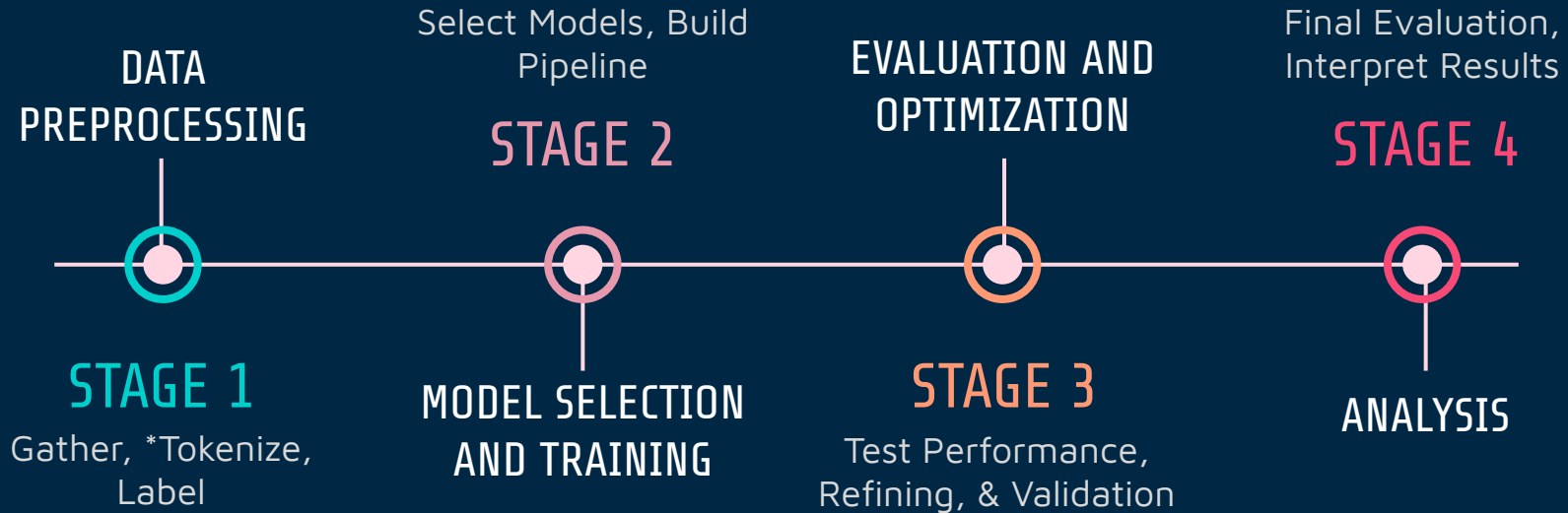
## Goals:

- ❖ If texts can be dated by **content** $\longrightarrow$ use ML to date **historical texts**?
- ❖ Can existing **machine learning** models help?
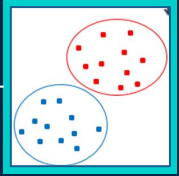
## Motivations:

- ❖ Interest in **Natural Language Processing (NLP)**.
    - ○ *Could we explore a problem that could utilize NLP & ML?*
- ❖ **Real-world application** of course content.

# APPROACH

**DATA PREPROCESSING**

**STAGE 1**

Gather, *Tokenize, Label

Select Models, Build Pipeline

**STAGE 2**

**MODEL SELECTION AND TRAINING**

**EVALUATION AND OPTIMIZATION**

**STAGE 3**

Test Performance, Refining, & Validation

Final Evaluation, Interpret Results

**STAGE 4**

**ANALYSIS**

**\*Tokenization**: Systematically breaking down a text into units (often called words, or subwords), and in the process, removing irrelevant features like capitalization, whitespace, punctuation, etc.
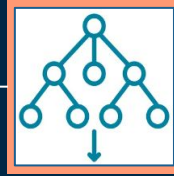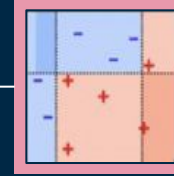
# MODELS



## 01

### NAIVE BAYES

1. Classifies extremely **quickly**.
2. Views all features as being **independant**.



## 02

### RANDOM FOREST

1. Can handle **high dimensional data** well.
2. Reduces **overfitting**.
3. Robust against **noise**.



## 03

### ADABOOST

1. Reduces **overfitting**.
2. Robust against **noise**.
3. Handles **class imbalance** well

# VECTORIZATION

**Word Vector:** A vector such that <u>words</u> are *mapped* to <u>indices</u>.

**Vectorization Algorithms:** Assign a <u>numeric feature value</u> per <u>word</u> *w* in the <u>word vector.</u>

*Example using BoW:*
  - **Ex. 1** ("The dog sat, the cat sat too.")
  - **Ex. 2** ("The dog and the cat and the other cat all sat.")

| Text # | $f_0$ 'the' | $f_1$ 'dog' | $f_2$ 'cat' | $f_3$ 'sat' |
|--------|-------------|-------------|-------------|-------------|
| Ex. 1  | 2           | 1           | 1           | 1           |
| Ex. 2  | 3           | 1           | 2           | 1           |

# VECTORIZATION ALGORITHMS

**Bag-of-Words (BoW):** "The <u>frequency</u> of <u>occurrence</u> of each word is used as a numeric feature for training a classifier."

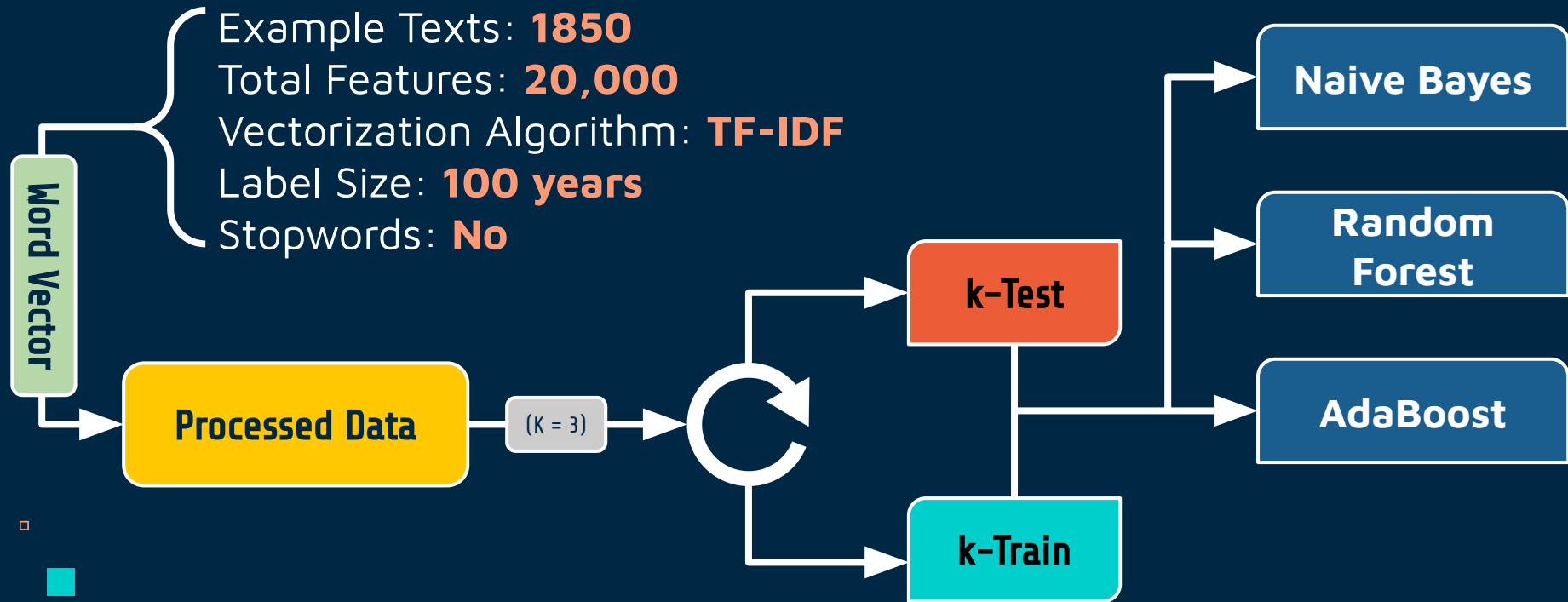**Term-Frequency Inverse Document Frequency (TF-IDF):**
- ❖ "Numerically reflects how <u>important</u> a word is to a document."
- ❖ TF-IDF value increases <u>proportionally</u> to the number of times a word appears in the document, and is <u>offset</u> by the number of documents in the corpus that contain the word.

$$W_{i,j} = tf_{i,j} \times \log \frac{(N)}{df_i}$$

# DATA GENERATION, TRAIN, & TEST PIPELINE

## Example Data Overview:

Example Texts: **1850**
Total Features: **20,000**
Vectorization Algorithm: **TF-IDF**
Label Size: **100 years**
Stopwords: **No**

Word Vector

Processed Data

(K = 3)

k-Test

k-Train

Naive Bayes

Random Forest

AdaBoost

# VALIDATION

**Randomized Search** with **cross validation** across a <u>range</u> of values for each of the values:
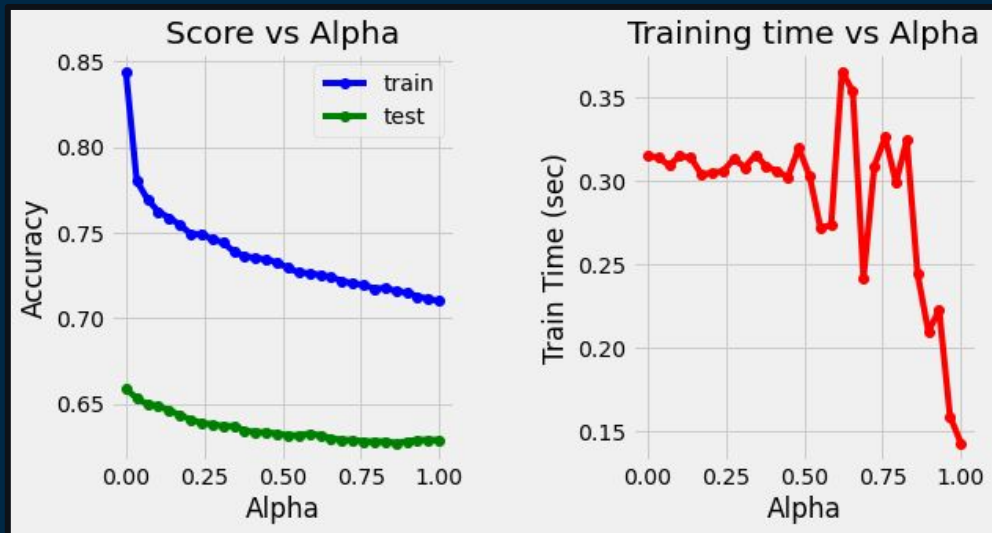
- ❖ **Naive Bayes** (Hyperparameters):
    - ■ *alpha* (0, 1]

- ❖ **Random Forests** (Hyperparameters):
    - ■ *n_estimators* [200, # unique words * 1/3]
    - ■ *min_samples_split* [2, 10]
    - ■ *min _samples_leaf* [1, 4]
    - ■ *max_features* ['auto', 'sqrt']
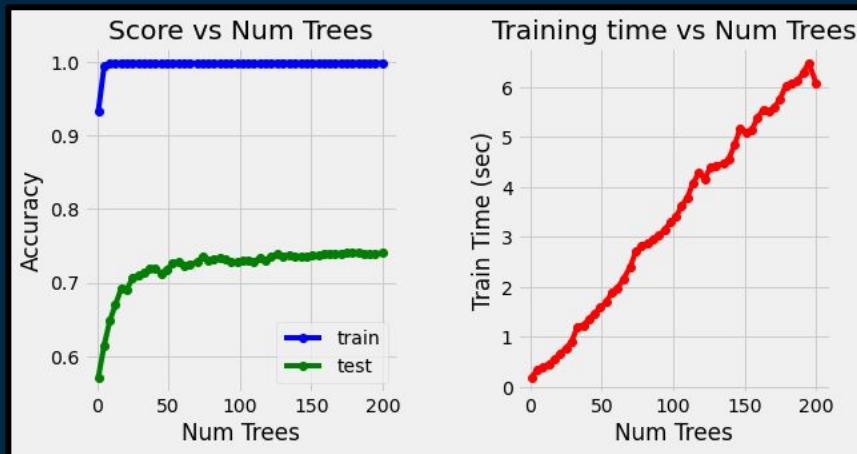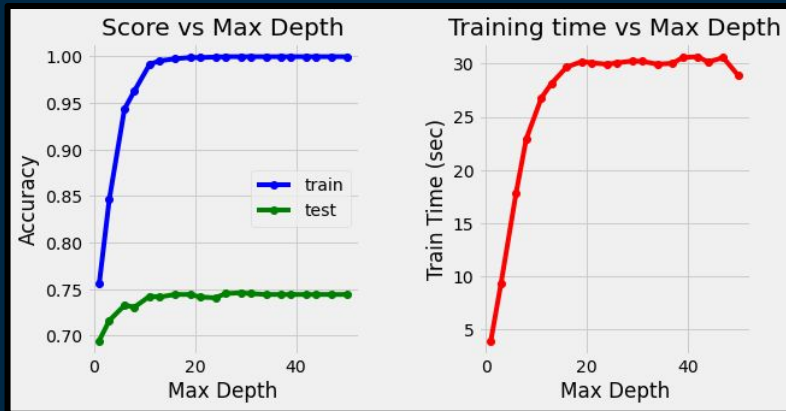    - ■ *max_depth* `[10, 110]
    - ■ *bootstrap* [true, false]

- ❖ **AdaBoost** (Hyperparameters):
    - ■ *n_estimators* [50, 500]
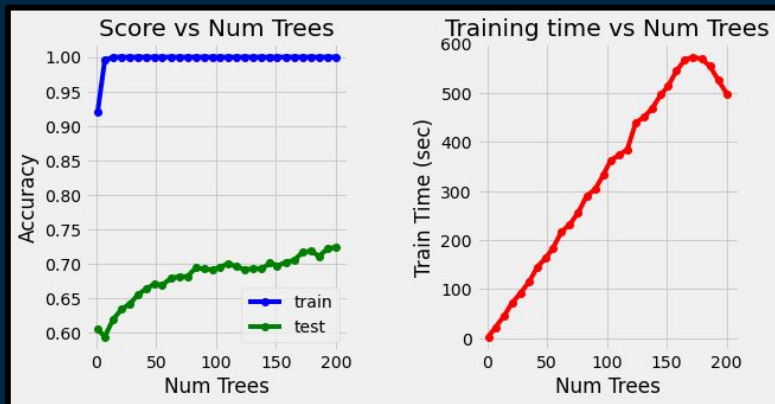    - ■ *learning_rate* [0.001, 1.0]
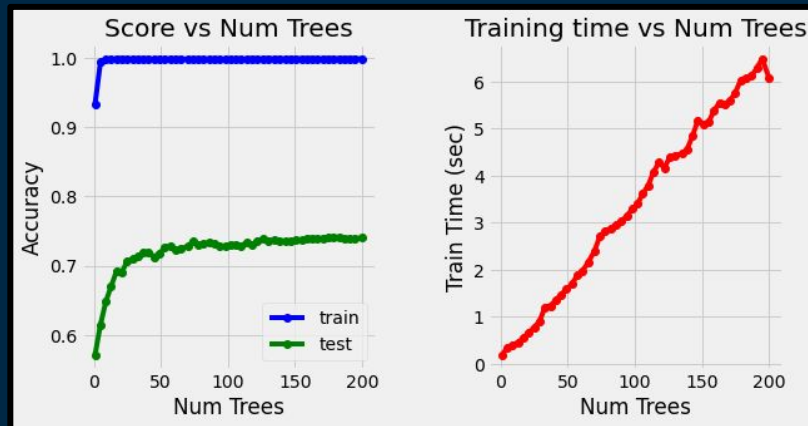    - ■ *base_estimator* [1, 15]
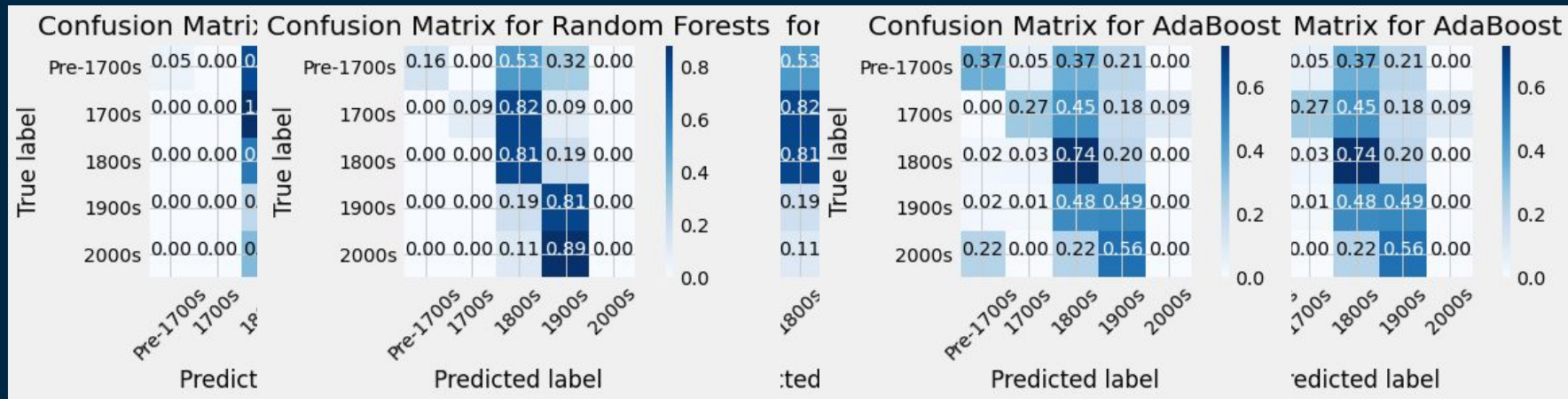
# NAIVE BAYES

# RANDOM FOREST

# ADABOOST

# RANDOM FOREST

# RESULTS

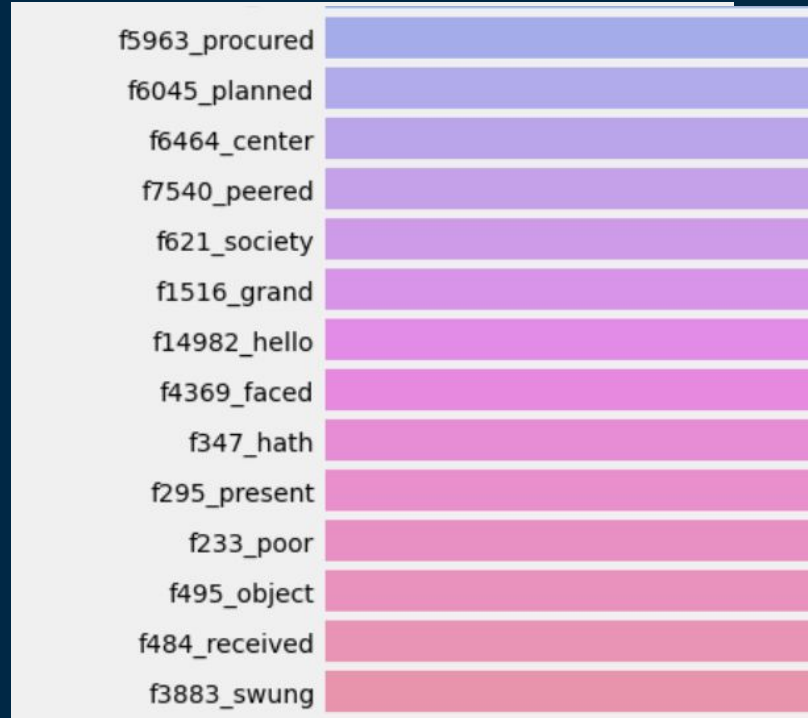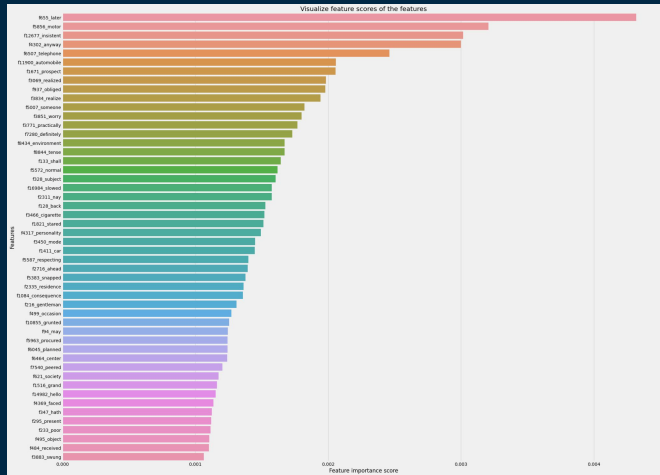| Model | Base | Tuned | Performance Delta |
|-------|------|-------|-------------------|
| AdaBoost | 0.5189 | 0.7270 | 0.2081 |
| Naive Bayes | 0.6527 | 0.7054 | 0.0527 |
| Random Forest | 0.7568 | 0.7773 | 0.0205 |

- Random Forest: Best Performing
  - We see the smallest gain in parameter tuning
- AdaBoost: Second Best (Most Potential)
  - Potentially more gains from more trees
- Naive Bayes: Fastest + Simplest
  - Simpler model (probably no further gains) but very fast.
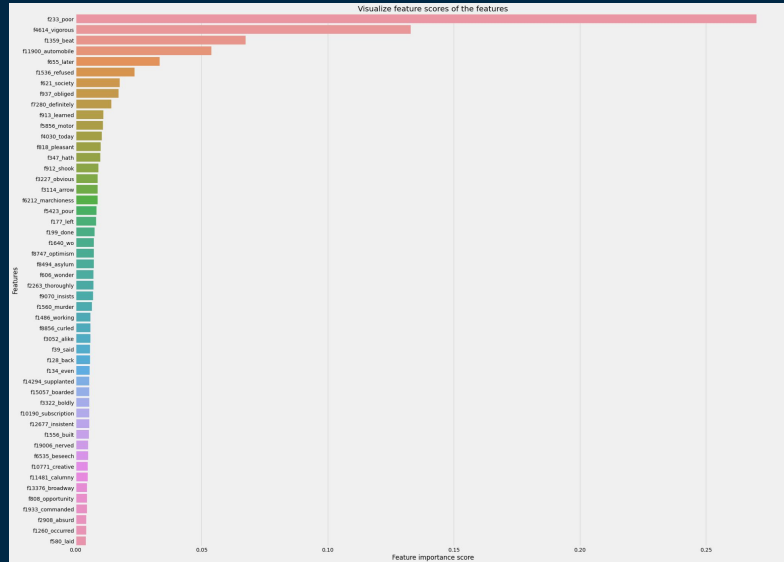
# RESULTS (cont.)

# Features

# Random Forests

# Features
# AdaBoost

# CRITIQUES, CHALLENGES & REFLECTIONS

- **Unbalanced Dataset (Classes)**
  - Date imbalance (severe class imbalance)
  - Unverifiable examples
- **Feature Standardization Approach**
  - ~100 texts arbitrary, could be missing important features
- **Boundary Classification**
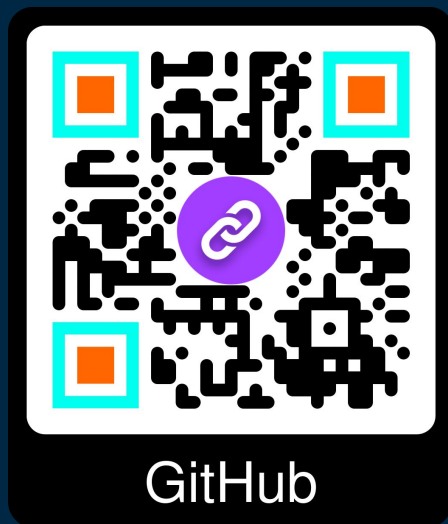  - Is there really a way to tell texts written in 1899/1900 apart?

# SUMMARY & FURTHER WORK

- **Can Texts Be Dated Purely By Content?**
  - Yes!
- **Can Existing Machine Learning Models Help?**
  - Yes!
- **Are They Accurate?**
  - ~70%!
- **Could We Apply This To Date Historical Texts?**
  - A question for further applied research.

## Further Work:

Feature Embedding; Sample Rebalancing;
Deep Learning; Feature Semantic Analysis

# REFERENCES


GitHub

1. *Project gutenberg*. Project Gutenberg. (1977). https://www.gutenberg.org/
2. Liebeskind, C., & Liebeskind, S. (2020). Deep learning for period classification of Historical Hebrew texts. *Journal of Data Mining & Digital Humanities, 2020*. https://doi.org/10.46298/jdmdh.5864
3. Gaudin, R. et al. (2014, July 6). *Openzim/Gutenberg: Scraper for downloading the entire ebooks repository of project gutenberg*. GitHub. https://github.com/openzim/gutenberg