

# Final Project Report

Group 3: Jehan Bugli

## Overview

This project tags congressional bills with one of [32 policy area designations](#), replacing an important yet tedious task that is currently completed manually by the Congressional Research Service (CRS).

The CRS assists legislators with **policy analysis**, **official bill summaries**, and more; however, their job is rapidly becoming more difficult due to **rapidly increasing bill introductions** causing significant backlogs for official text releases, which in turn affects third parties reliant upon their processing efforts.

This experiment's goal is to provide practical utility; as a result, this primarily leverages classical methods. These methods have two key advantages for policy area tagging relative to state-of-the-art neural networks:

- **Interpretability:** results can be traced back with relatively straightforward reasoning, which is critical for real-world adoption; analysts should be able to explain the logic behind these labels to various stakeholders (members of Congress, lobbyists, etc.)
- **Cost:** models are significantly cheaper to train and run for inference

This report will:

- Discuss the data collection and processing efforts involved
- Describe the chosen modeling approach and experiment setup
- Outline the experiment results and conclusions

## Data collection and processing

Data was sourced from the GovInfo bulk data repository which stores XML-formatted bill data for recent congressional sessions; this dates back to the 113<sup>th</sup> session (2013-2015).

To start, I downloaded bill text for each bill in the available sessions. This includes every version; for instance, if a bill received amendments between introduction and enactment, both the amended and original versions would be included. These were saved as a collection of session-specific .zip files.

**Commented [j1]:** <https://www.congress.gov/crs-products>

**Commented [j2]:** <https://www.congress.gov/help/bill-summaries>

**Commented [j3]:** <https://rollcall.com/2025/03/05/publishing-pileup-congressional-bills-slow-to-reach-public/>

**Commented [j4]:** <https://www.govinfo.gov/bulkdata>

I also downloaded associated data for each bill, including policy area tags and other CRS-processed items, in the same format as noted earlier.

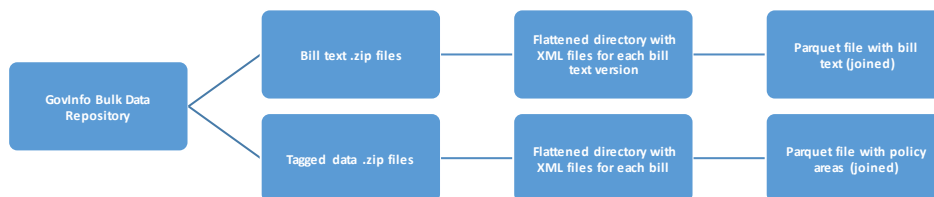
Before further processing, I opted to “flatten” both folders, unzipping contents into two main directories.

Then, I converted each flattened folder into a tabular dataset. The folder with bill text was processed to generate a parquet file storing each file’s name and bill text.

Using BeautifulSoup’s XML parser, I searched for valid legislative body tags (e.g. ‘legis-body’, ‘preamble’) and discarded components such as the title, publishing date, sponsors, sessions, and more. While some of these items certainly may have provided some predictive value, I was concerned that such modeling may step outside of the NLP realm and muddy the results’ interpretability.

The folder with additional bill data was processed to generate a parquet file storing policy areas alongside text version names.

These text version names were used to join the two files and construct my finalized input dataset, containing 115,100 tagged bill text samples.



## Modeling approaches

At its core, my approach uses Term Frequency – Inverse Document Frequency (TF-IDF). TF-IDF emerged out of pioneering information retrieval research as a method to quantify word importance.

This combines two key metrics:

- **Term Frequency (TF):** How frequently a term appears within a single document
- **Inverse Document Frequency (IDF):** How rare a term is across the entire document collection

**Commented [j5]:** <https://www.byteplus.com/en/topic/400324?title=the-origins-of-tf-idf-term-frequency-inverse-document-frequency>

The TF-IDF score multiplies these metrics to weight terms that are frequent within a given document but relatively rare across the entire collection of documents. This weighting can prove useful for many tasks, including feature extraction for classification.

After testing parameter combinations, I tuned TF-IDF vectorization by:

- Removing common English stop-words
- Capturing individual words and meaningful two-word phrases
- Ignoring terms that only appear in a single document
- Removing terms that appear in over 85% of documents

I trained three different models on these extracted features:

- **LinearSVC (Support Vector Classifier):** attempts to find the "best" separating hyperplane that divides the TF-IDF feature space, distinguishing between different policy area categories. It aims to maximize the distance between the closest data points (support vectors) of different classes.
- **Logistic Regression:** predicts the probability of an instance belonging to a particular class using a logistic (sigmoid) function applied to a linear combination of the input features (TF-IDF scores).
- **Multinomial Naive Bayes:** calculates the probability of a document belonging to each policy area based on the frequency of terms within it, making a "naive" assumption that the features are independent of each other given the class.

**Commented [j6]:** <https://scikit-learn.org/stable/modules/svm.html#svm-classification>

**Commented [j7]:** [https://scikit-learn.org/stable/modules/linear\\_model.html#logistic-regression](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression)

**Commented [j8]:** [https://scikit-learn.org/stable/modules/naive\\_bayes.html#multinomial-naive-bayes](https://scikit-learn.org/stable/modules/naive_bayes.html#multinomial-naive-bayes)

I compared the options above using grid search cross-validation; this explores various parameter combinations across the TF-IDF vectorizer and these classifiers, comparing performance on weighted recall to return the best-performing combination.

For each combination, this search performs 3-fold cross-validation, where the data is split into 3 subsets; for each "fold", a model is trained using 2 subsets and the remaining one serves as a validation set. This prevents overfitting, where the model is unable to generalize effectively from its training data.

**Commented [j9]:** [https://scikit-learn.org/stable/modules/cross\\_validation.html](https://scikit-learn.org/stable/modules/cross_validation.html)

My grid search included different hyperparameter configurations, such as:

- **C:** A regularization parameter penalizing misclassified examples; a larger C value prioritizes accurate classification but risks overfitting.
- **Maximum iterations:** The number of allowed iterations, altered to ensure that models are converging

I focused on macro recall as my primary metric. Recall is the ratio of true positives over all actual positive instances for a class (including both true positives and false negatives); my methodology averages recall across all classes without regard for imbalance. This decision is made out of perceived practical utility; for a policy affairs professional, catching every relevant piece of legislation (each “true positive”) for their potentially niche field is of the utmost importance, even if that arrives alongside a number of false positives.

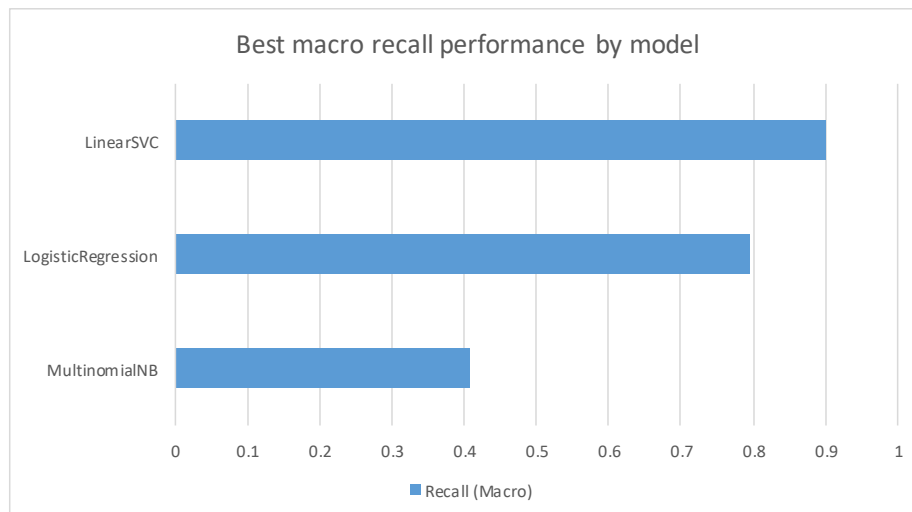
## Results

The linear SVC model displayed the strongest performance by far, reaching ~90% macro recall with grid search hyperparameter tuning!

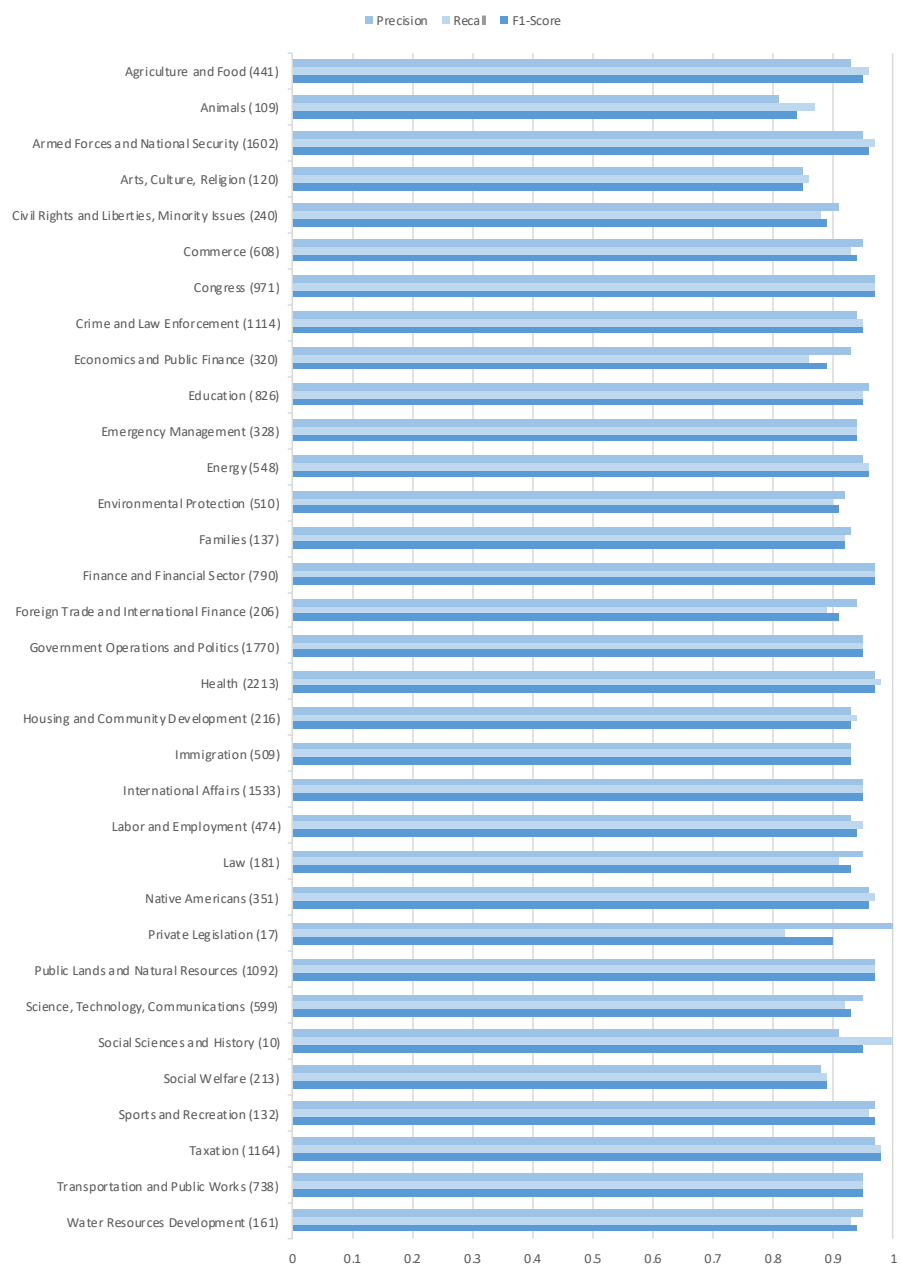
This served as an encouraging sign for the classical approach. Logistic regression followed close behind, while multinomial naïve Bayes was a distant third.

The graphs below display:

- Macro recall performance across the three tested models, including their best results from grid search tuning
- Area-by-area recall, precision, and F1-score results for the best-performing LinearSVC model



## LinearSVC precision, recall, and F1 score on the test set



## Discussion

These results support the notion that a classical approach can form a very robust classification mechanism for the federal legislation corpus.

While the LinearSVC model displayed excellent results, however, performance differed across different policy areas as displayed in the prior graph!

For instance, “Private Legislation” saw especially low recall values, which is reasonable given its relative infrequency and lack of topic specificity. There was stability with respect to class imbalance, likely thanks to the macro recall objective; “Social Sciences and History”, with only 10 test set instances, saw perfect recall.

While discrepancies are noteworthy and warrant further investigation, I think that the generally robust performance is a positive sign for practical implementation. All class predictions carry a reasonably high degree of certainty, and while policy areas are assigned in a mutually exclusive manner, the bills themselves are not as neatly partitioned! For instance, “Foreign Trade and International Finance” could be viewed as adjacent to “Finance” as well as “International Affairs”.

The poor MultinomialNB performance suggests that the “naïve” assumption doesn’t hold, such that individual term impacts on categorization cannot be viewed independently!

This makes sense intuitively; compound terms or phrases found jointly in a document could imply classifications that the individual terms would not, such as “pharmacy benefit manager”.

I think that results could be improved to a degree by refining stop-words to target common legal phrasing; the current results use the standard TfidfVectorizer English set, which may be a poor fit.

## Conclusion

Overall, the results ended up as a welcome surprise; while I was initially concerned about a potential trade-off between performance and interpretability, classical methods appear fully capable of supporting real-world tagging applications.

There are a number of avenues for further investigation and application.

As noted earlier, while policy areas are mutually exclusive, legislative contents can include relevant impacts on multiple. Simply judging this model based on the assigned class may not be

**Commented [j10]:** <https://libguides.law.umn.edu/c.php?g=125795&p=823607>

sufficient to judge practical use; building out a representation of “adjacent” policy areas and their relationships could allow us to evaluate model performance in a more realistic manner.

Testing whether this approach could be generalized to other text corpora is also worth considering. For instance, state-level legislation, Federal Register rules, and others are classified fully independently (if at all)! A unified classification system could help the public more easily navigate the mass of available legal materials to identify items that are relevant to them.

This also carries potential for retroactive or custom classification; while legislative topics have shifted over time, the labels across the digitized corpus remain static. If the CRS wants to add a new label (e.g. “Artificial Intelligence”), they could potentially re-train a model on an expanded label set (adding the new label manually for a session) and retroactively reclassify older bills to reflect that change. Alternatively, if a public stakeholder (such as a policy affairs team) has an in-house set of labels that they apply to bills, this modeling approach could be used to improve tracking mechanisms.

Finally, I think that this method has value for investigating legislative attitudes at large across the training corpus. For instance, while the current model associates “China” and “Africa” with “International Affairs”, it associates “Chinese” and “African” with “Armed Forces and National Security”. Understanding such term interactions may reveal subtle trends and biases that a less interpretable modeling approach might not reveal.