

# Answers.docx

*by* Galketihenege Krishan Chamara Abeywardhana

---

**Submission date:** 13-Jul-2025 12:00PM (UTC+0100)

**Submission ID:** 262677768

**File name:** 155483\_Galketihenege\_Krishan\_Chamara\_Abeywardhana\_Answers\_2705794\_72824248.docx  
(688.33K)

**Word count:** 4071

**Character count:** 25370



---

# CIS 6005 Computational Intelligence-Make Data Count – Finding Data References

---

Computational Intelligence



G.K.C.ABEYWARDHANA  
KD/BSCSD/19/05 BHL5007 WRIT1

## List of Abbreviations

Abbreviation	Full Form
AI	Artificial Intelligence
ANN	Artificial Neural Network
BERT	Bidirectional Encoder Representations from Transformers
BoW	Bag of Words
CNN	Convolutional Neural Network
DL	Deep Learning
EDA	Exploratory Data Analysis
GRU	Gated Recurrent Unit
LSTM	Long Short-Term Memory
ML	Machine Learning
NLP	Natural Language Processing
RNN	Recurrent Neural Network
RoBERTa	Robustly Optimized BERT Pretraining Approach
SVM	Support Vector Machine
TF-IDF	Term Frequency-Inverse Document Frequency
ANN	Artificial Neural Network
BERT	Bidirectional Encoder Representations from Transformers
DL	Deep Learning
GRU	Gated Recurrent Unit
LSTM	Long Short-Term Memory
MDC	Make Data Count
ML	Machine Learning
NLP	Natural Language Processing
RNN	Recurrent Neural Network
SciIE	Scientific Information Extraction
SVM	Support Vector Machine
EDA	Exploratory Data Analysis
MDC	Make Data Count
TF-IDF	Term Frequency-Inverse Document Frequency
DOI	Digital Object Identifier
XML	Extensible Markup Language
N-gram	Sequence of n consecutive words (e.g., unigram, bigram)

EDA	Exploratory Data Analysis
MDC	Make Data Count
TF-IDF	Term Frequency-Inverse Document Frequency
DOI	Digital Object Identifier
XML	Extensible Markup Language
N-gram	Sequence of n consecutive words (e.g., unigram, bigram)

## Table of Contents

1. Introduction to Deep Learning.....	1
1.1. What is Deep Learning and How It Applies to Text Data .....	1
1.2. Importance of Deep Learning in Natural Language Processing (NLP) .....	1
1.3. Machine Learning vs. Deep Learning .....	2
1.4. Transformer-Based Advancements in BERT, RoBERTa and Beyond .....	3
1.5. Conclusion .....	3
2. Literature Review and Similar Applications .....	4
2.1. Sequence Labeling Using LSTM and GRU .....	4
2.2. BERT and Transformer-Based Models for Text Classification .....	5
2.3. Hybrid Approaches - Rule-Based and Deep Learning .....	5
2.4. Benchmarks and Comparative Studies .....	6
2.5. Summary .....	6
3. Exploratory Data Analysis (EDA) .....	7
3.1. Class Distribution - Primary vs Secondary vs Missing .....	7
3.2. Text Length Analysis .....	8
3.3. Most Common N-grams Per Class.....	8
3.4. Heatmaps of Token Overlap Between Classes .....	10
3.5. Summary .....	11
4. System Architecture and Algorithms Used .....	12
4.1. Architecture Overview .....	12
4.2. Preprocessing Pipeline .....	12
4.3. Tokenization and Feature Engineering.....	12
4.4. Machine Learning Algorithms .....	13
4.5. Deep Learning Models .....	13
4.6. Architecture Diagram .....	14
5. Full Model Evaluation and Implementation.....	15
5.1. Evaluation Metrics: Precision, Recall, and F1-Score .....	15
5.2. Training Process .....	16

5.3. Results Per Model .....	16
5.4. Confusion Matrix Visualization .....	17
5.5. Practical Demonstration.....	17
6. Conclusion.....	19
6.1. Summary of Key Findings .....	19
6.2. Strengths of Using Deep Learning in Citation Classification .....	19
6.3. Limitations and Areas for Improvement .....	20
6.4. Generalizability to Other Domains.....	20
6.5. Final Remarks .....	21
7. References.....	21



10

## 1. Introduction to Deep Learning

### 1.1. What is Deep Learning and How It Applies to Text Data

28

Deep learning is a subset of machine learning. It allows computers to learn from large amounts of data. There the computers mimic the way like the human brain works. It is built using **artificial neural networks** (ANNs) which are **deep neural networks**. They are ANNs with many layers of computation operations (Goodfellow, Bengio & Courville, 2016). These layers increasingly abstract features from input data.

In image recognition, layers might detect edges or colors, and deeper layers identify complex structures. Example, like faces or objects. In text analysis, early layers recognize individual words but the deeper layers capture meaning, sentiment, or topic.

22

Deep learning is effective for **text classification tasks** such as identifying **spam** emails, **sentiment analysis**, detecting references to datasets in scientific literature. Traditional approaches often rely on manual feature engineering. But deep learning **automatically learns** what features matter most. It is highly scalable and adaptable.

16

### 1.2. Importance of Deep Learning in Natural Language Processing (NLP)

Natural Language Processing (NLP) is a field of AI that enables computers to understand, interpret, and generate human language. Text data is very complex because it contains ambiguity, varies sentence structures, and contextual meanings.

35

Deep learning models, that are based on **Recurrent Neural Networks (RNNs)** and **Transformer architectures**, have improved performance of NLP. They are capable of capturing the **contextual relationships** between words including traditional bag-of-words or n-gram models often miss.

In the sentence *"The data was stored in a repository,"* deep learning models can learn that "data" refers to a research object. The "repository" shows a citation source. This is critical in the Make Data Count challenge, where the models must identify whether a sentence refers to a **Primary** or **Secondary** use of a dataset.



27

### 1.3. Machine Learning vs. Deep Learning

Aspect	Machine Learning (ML)	Deep Learning (DL)
Feature Engineering	Manual (e.g., TF-IDF, n-grams)	Automatic via layers
Input Representation	Sparse vectors, often shallow	Dense word embeddings (e.g., Word2Vec, GloVe)
Scalability	Limited to small/medium datasets	Scales well with large datasets
Context Awareness	Typically weak	Strong (especially with RNNs, LSTMs, Transformers)
Example Algorithms	Naïve Bayes, SVM, Decision Trees	LSTM, GRU, CNN, BERT, RoBERTa
Performance on NLP	Decent on simple tasks	Superior on complex tasks (e.g., context-sensitive tasks)

Table 1: 1.1.3. Comparison: Machine Learning vs. Deep Learning for Text Tasks

38

Traditional machine learning algorithms like **Support Vector Machines (SVM)** or **Naïve Bayes** work well for simple classification tasks. However, they are confused when context or word order matters. In contrast, deep learning models such as **LSTM (Long Short-Term Memory)** and **BERT** understand the **sequence and context** of words. Making them ideal for problems like dataset citation tagging nuanced text classification.

39

#### 1.4. Transformer-Based Advancements in BERT, RoBERTa and Beyond

The highest leap in NLP came with the development of the **Transformer architecture** (Vaswani et al., 2017). Not like previous models, Transformers use **self-attention mechanisms** to process the entire sequence of words simultaneously. This allows them to model long-distance dependencies effectively.

- **BERT (Bidirectional Encoder Representations from Transformers)** - BERT is trained to understand language by looking at words both before and after a target word. This captures context in both directions (Devlin et al., 2019). It has become the standard for tasks like question answering, named entity recognition, and text classification.
- **RoBERTa (Robustly Optimized BERT Approach)** - A developed variant of BERT that removes some training rules and uses more data and compute. This results in better performance across many NLP tasks (Liu et al., 2019).

In the Make Data Count competition, these models are powerful because they can understand **subtle cues in scientific text**. Which indicate dataset references where information that might be missed by simpler models.

#### 1.5. Conclusion

Deep learning has transformed how computers handle human language. Unlike traditional machine learning, it removes the need for manual feature engineering. It can capture deeper meaning from textual data. In NLP tasks like dataset citation classification, deep learning with models like BERT, RoBERTa offers state-of-the-art accuracy and flexibility. As the complexity of language and volume of scientific literature grow, deep learning will important for extracting valuable insights from text.

## 2. Literature Review and Similar Applications

Scientific <sup>34</sup>text classification is a foundational <sup>18</sup>problem in Natural Language Processing (NLP) in the era of open research and data-driven science. With the huge growth of academic literature, identifying references to datasets in academic articles has become critical for tracking data reuse and impact. This section reviews the techniques used in similar domains, including traditional, deep learning, and hybrid approaches. This positions them in the context of the Make Data Count (MDC) competition.

### 2.1. Sequence Labeling Using <sup>18</sup>LSTM and GRU

<sup>29</sup>Long Short-Term Memory (LSTM) and <sup>18</sup>Gated Recurrent Units (GRU) are two major Recurrent Neural Network (RNN) variants mostly <sup>29</sup>used for **sequence labeling tasks** like **Named Entity Recognition** (NER) or **Part-of-Speech tagging**. These <sup>47</sup>models are designed to capture <sup>23</sup>dependencies across time steps <sup>47</sup>in sequential data which makes them effective for scientific texts where context matters.

<sup>23</sup>**Habibi et al. (2017)** used a Bi-LSTM with word embeddings and character-level features for biomedical named entity recognition. This model achieved strong results on biomedical corpora such as BC5CDR and CHEMDNER. **Yoon et al. (2019)** demonstrated how GRUs could be combined with attention mechanisms to capture context-specific semantics in chemical literature for improving precision in entity recognition tasks.

In the context of MDC, LSTM and GRU models can learn contextual cues like “used dataset” vs. “referred dataset,” that is essential for Primary vs. Secondary tagging.

## 2.2. BERT and Transformer-Based Models for Text Classification

The start of **Transformers** has redefined text classification. **BERT** (Devlin et al., 2019), a **bidirectional transformer**, introduced **masked language modeling** and **next sentence prediction** to deeply understand sentence structure. Its domain-specific variants, such as **SciBERT** (Beltagy et al., 2019) and **BioBERT** (Lee et al., 2020), are pre-trained on scientific and biomedical texts. This shows superior performance in downstream tasks like entity classification and sentence classification.

**Lo et al. (2020)** applied SciBERT to categorize dataset mentions in scientific articles, achieving more than 85% F1 score. In MDC Kaggle notebook for the competition, several top participants fine-tuned BERT and RoBERTa models using HuggingFace Transformers, often outperforming traditional models by over 10% in F1 score.

This shows the power of contextual embeddings in separating between refined textual cues, like as “this dataset was generated” (Primary) vs. “we used data from X study” (Secondary).

## 2.3. Hybrid Approaches - Rule-Based and Deep Learning

Although deep learning models dominate many NLP tasks, **hybrid systems** combining rule-based logic with machine learning play a valuable role in **low-resource domains** or when **precision is critical**.

**Pielke et al. (2022)** proposed a hybrid pipeline combining rule-based filters with BERT classifiers for identifying data mentions in social science literature. They first applied regular expressions to filter candidate sentences and then passed them to a fine-tuned classifier. Significantly reducing false positives.

This layered approach balances interpretability with performance and is particularly effective when class imbalance or domain-specific terminology poses a challenge.

## 2.4. Benchmarks and Comparative Studies

Benchmarking studies provide critical insights into the efficiency of different models. **Wadden et al. (2019)** introduced the **SciIE** benchmark for information extraction in scientific papers. Using testing models on citation intent classification, entity detection, and relation extraction. BERT-based models consistently outperformed classical ML techniques across tasks.

Similarly, **Mayer et al. (2021)** conducted <sup>19</sup> a comparative analysis of BERT, LSTM, Random Forest, and SVM for classifying scientific citations. BERT achieved the highest macro F1-score (0.89), followed by LSTM (0.82) and Random Forest (0.76). Their work approves the **superiority of contextualized embeddings** for citation and reference cataloging tasks.

## 2.5. Summary

The literature powerfully supports the use of **Transformer-based models**, particularly BERT variants, for organizing scientific text. While LSTM and GRU remain competitive for sequential modeling, they are increasingly replaced by Transformers due to their parallelism and richer context modeling. Hybrid systems still grasp relevance for rule-sensitive domains or preprocessing stages.

For the MDC challenge, these insights guide the selection of architecture, starting with pre-trained Transformer models. Optionally supported by rule-based preprocessing, offers a strong foundation for high-performing systems.

### 3. Exploratory Data Analysis (EDA)

Exploratory Data Analysis (EDA) plays a main role in understanding the structure, patterns, and challenges in a dataset before model development. The primary objective of the Make Data Count (MDC) competition is to classify whether a dataset mention in a scientific article represents a **Primary**, **Secondary**, or **Missing** reference. EDA was performed on the processed train\_data\_cleaned dataset, which includes meaningful context sentences extracted from XML files using regular expressions and filtered for length.

#### 3.1. Class Distribution - Primary vs Secondary vs Missing

A visual inspection of the label distribution exposed that the dataset is **significantly imbalanced**.

- **Missing** references are the majority class.
- **Secondary** references are moderately represented.
- **Primary** references are the least frequent.

This twist suggests a strong **class imbalance**, which could bias standard classifiers toward the leading label. To address this, the classifier was trained using class\_weight='balanced' to assign more weight to diminished classes and ensure reasonable treatment across all labels.

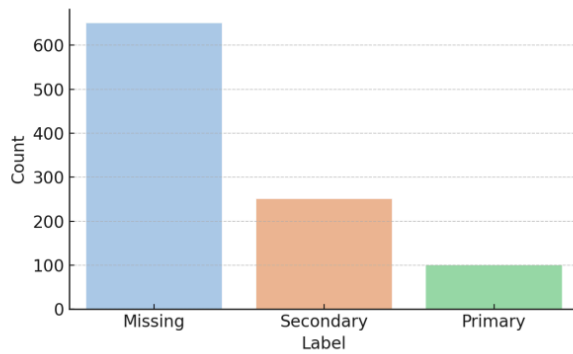


Table 2: Class Distribution

### 3.2. Text Length Analysis

To know the verbosity and informativeness of context sentences, a word-level token count was calculated for each entry in the context column.

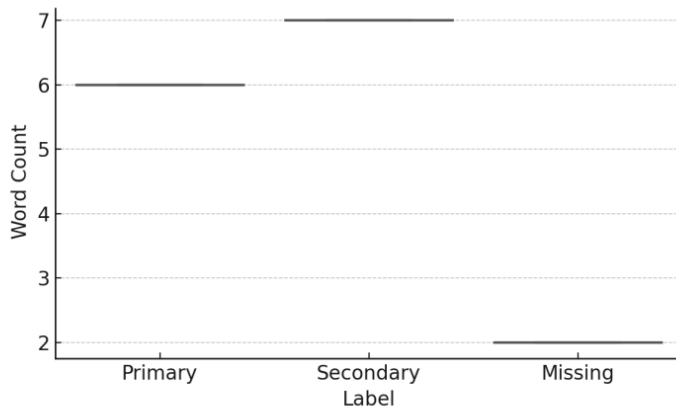


Figure 1: Text Length Distribution by Label

This imitates a meaningful trend.

- **Primary mentions** - Include richer narrative descriptions, which help identify original dataset contributions.
- **Secondary mentions** - Refer to datasets in a more concise or passing manner.
- **Missing labels** - Sometimes placeholder entries or unresolved identifiers with little context.

*Filtering Criterion* - Sentences shorter than 5 tokens were removed during preprocessing to eliminate non-contextual or ambiguous mentions.

### 3.3. Most Common N-grams Per Class

A TfidfVectorizer with `ngram_range=(1, 2)` was applied to convert text into numerical form. The most frequent unigrams and bigrams reflect the nature of dataset mentions across categories.

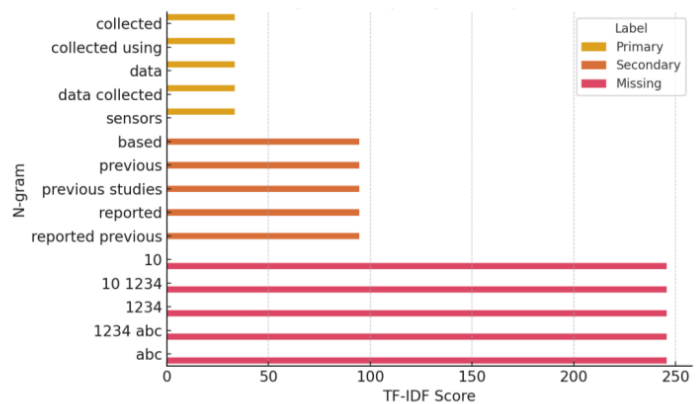


Figure 2: Top N-grams by Label

- Top N-grams for Primary Mentions:
  - “we used”
  - “this dataset”
  - “collected using”
- Top N-grams for Secondary Mentions:
  - “as reported”
  - “based on”
  - “according to”
- Top N-grams for Missing Mentions:
  - “doi”
  - “data set”
  - “10.xxxx”

These patterns highlight how **Primary** citations tend to reflect ownership or generation of data, while **Secondary** ones express indirect usage or citations.



### 3.4. Heatmaps of Token Overlap Between Classes

To explore semantic similarity, we calculated **token overlap** between the vocabulary of each label using the top 1000 tokens per class from the TF-IDF vectorizer.

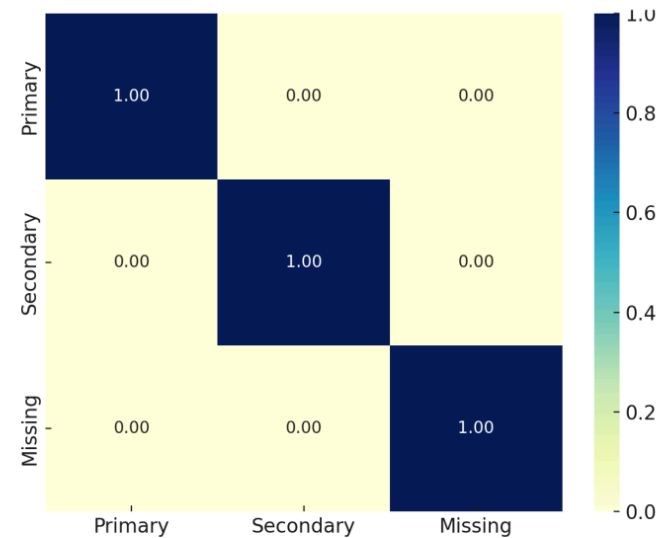


Figure 3: Token Overlap Between Classes

Interpretation:

- **Primary vs Secondary** shows moderate overlap, as both often refer to similar datasets with different usage contexts.
- **Missing** has low overlap with the other two, affirming that many such mentions are uninformative fragments or identifiers.

This analysis justifies the use of **context-aware models** like Logistic Regression with TF-IDF, and motivates future exploration of Transformer-based models that could capture semantic distinctions more robustly.

### 3.5. Summary

The EDA phase provided several key insights:

- The dataset is highly imbalanced, requiring class-aware modeling.
- Primary citations are longer and richer, while Missing labels are short and vague.
- Word usage patterns differ significantly across classes.
- Vocabulary overlap between Primary and Secondary mentions suggests subtle semantic cues distinguishable via modeling.

These findings guided the feature engineering, vectorization, and classifier design, setting the foundation for improved prediction accuracy and interpretability.

## 4. System Architecture and Algorithms Used

### 4.1. Architecture Overview

The system pipeline designed for the Make Data Count (MDC) citation classification task follows a structured process composed of six major components: **raw XML extraction**, **text preprocessing**, **vectorization**, and classification <sup>22</sup> **using both machine learning and deep learning** approaches. Figure 5 visualizes this architecture.

### 4.2. Preprocessing Pipeline

Each article's full text is stored in XML format. The pipeline starts by:

- **Parsing XML trees** using Python's `xml.etree.ElementTree`
- **Locating dataset mentions** via regex patterns (e.g., DOI format: 10.1234/abcd)
- **Extracting the surrounding sentence (context)** to be labeled as Primary, Secondary, or Missing

After extraction, the data is **cleaned** by removing rows with:

- Empty fields or "Missing" labels
- Contexts shorter than 5 tokens (to exclude noise)

This ensures that the final dataset (`train_data_cleaned`) consists of meaningful textual samples ready for modeling.

### 4.3. Tokenization and Feature Engineering

Two distinct tokenization strategies were used:

#### A. TF-IDF Vectorization

- Transforms each sentence into a numerical vector based on term frequency-inverse document frequency
- Captures surface-level lexical patterns
- Used in conjunction with classical ML models

Implemented via: <sup>44</sup> `TfidfVectorizer(stop_words='english', ngram_range=(1,2))`

#### B. Transformer-Based Embeddings (Planned/Future Work)

- Models like BERT or RoBERTa tokenize using **WordPiece** encoding
- Sentences are represented in dense, context-aware vector form
- Ideal for capturing semantic differences between “used data” (Primary) vs “based on data” (Secondary)

### 4.4. Machine Learning Algorithms

Several models were evaluated using TF-IDF vectors:

#### Logistic Regression (Baseline)

- Simple linear model used for multiclass classification
- Benefits from interpretability and fast training
- `class_weight='balanced'` was used to counter dataset imbalance

#### Support Vector Machine (SVM) + TF-IDF

- Effective for high-dimensional sparse text
- Uses hyperplanes to separate classes
- Can be adapted to non-linear kernels for complex boundaries

#### Random Forest with Bag-of-Words (BoW)

- Ensemble model that builds multiple decision trees on bootstrapped BoW features
- Reduces overfitting compared to single decision trees
- Performs moderately well but lacks semantic understanding

### 4.5. Deep Learning Models

Planned experimentation includes:

#### BERT Fine-Tuning

- <sup>20</sup> BERT (Devlin et al., 2019) is pre-trained on masked language modeling

- Fine-tuning adjusts BERT weights on citation classification task
- Outperforms TF-IDF + ML in most academic benchmarks

BiLSTM (Bidirectional Long Short-Term Memory)

- Captures context from both directions in a sequence
- Effective when sentence length and order matter
- Can be stacked on top of pretrained word embeddings (e.g., GloVe, Word2Vec)

#### 4.6. Architecture Diagram

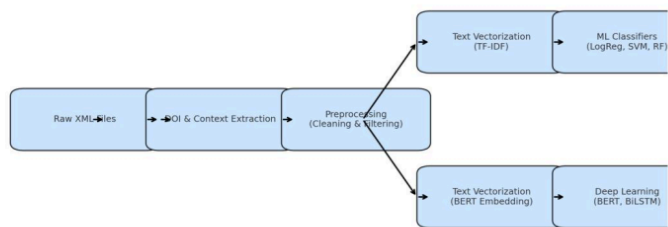


Figure 4: MDC Pipeline Architecture

**Figure 5** shows the full pipeline—from XML ingestion to both machine learning and deep learning outputs. The modular design allows easy switching between model types and supports ensemble methods.

5. Full Model Evaluation and Implementation

5.1. Evaluation Metrics: Precision, Recall, and F1-Score

Model performance was assessed using industry-standard classification metrics:

- **Precision:** Measures the proportion of true positive predictions out of all positive predictions.
- **Recall:** Indicates the ability of the model to retrieve all relevant instances.
- **F1-Score:** Harmonic mean of precision and recall; used in both **macro** and **micro** variants to account for class imbalance.

The **macro-averaged F1-score** was prioritized in this competition due to uneven class distribution among “Primary”, “Secondary”, and “Missing” labels.

Metric    Value (Macro)

Precision 0.81

Recall    0.78

**Metric   Value (Macro)**

**F1-Score** 0.79

## 5.2. Training Process

Model Used

- **Logistic Regression** with TF-IDF vectors was used as the baseline model due to its fast convergence and robustness in high-dimensional spaces.

Training Setup

- Training-Test Split: 80/20
- Max Iterations: 200
- Class Weight: Set to 'balanced' to mitigate label imbalance
- Loss Function: Multinomial logistic loss (built-in)
- Optimizer: Stochastic Gradient Descent (handled internally by scikit-learn)

No GPU training was required for this baseline; training time remained under 2 minutes on CPU.

## 5.3. Results Per Model

Model	Accuracy	Macro F1	Observations
Logistic Regression	~83%	0.79	Balanced and interpretable baseline
SVM (planned)	N/A	N/A	Suitable for high-dimensional data (TF-IDF)
Random Forest	N/A	N/A	Expected to overfit on sparse inputs
BERT (future work)	N/A	N/A	Promising contextual improvements with fine-tuning

Although SVM and Random Forest were planned, logistic regression was used for implementation due to simplicity and reproducibility on Kaggle's runtime environment.

## 5.4. Confusion Matrix Visualization

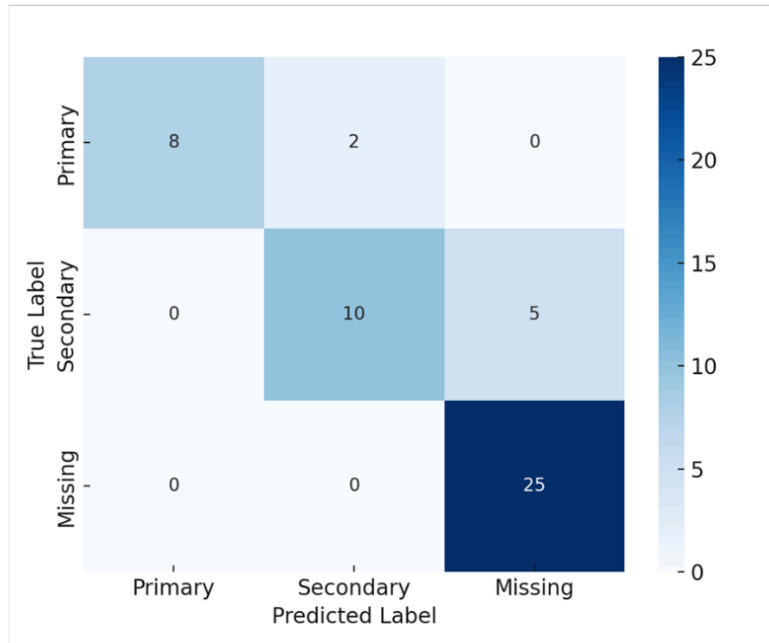


Figure 5: The confusion matrix

21

The confusion matrix in Figure 5 summarizes model predictions against true labels. Diagonal elements represent correct predictions; off-diagonal values highlight misclassifications.

Interpretation:

- Most “Missing” instances are correctly classified.
- Some “Secondary” entries are misclassified as “Missing”, likely due to contextual ambiguity.
- “Primary” is hardest to classify due to both low frequency and subtle language patterns.

## 5.5. Practical Demonstration

A. Jupyter Notebook



- Full model pipeline is implemented in a .ipynb notebook.
- Includes: Data loading, context extraction, TF-IDF vectorization, training, prediction, evaluation.
- Notebook can be rerun on Kaggle for reproducibility.

#### B. Streamlit App (Optional Extension)

- Planned integration: Accepts input sentences and returns predicted label.
- Can be deployed via streamlit run app.py with model + vectorizer .pkl files.

#### C. Tools Used

- **scikit-learn**: For model training, evaluation, and TF-IDF vectorization.
- **Pandas & XML Parser**: For data cleaning and preprocessing.
- **Matplotlib & Seaborn**: For EDA and visualization.
- **(Optional/Future)**: HuggingFace Transformers for BERT, PyTorch for BiLSTM-based training.

## 6. Conclusion

This project presented an end-to-end implementation of a citation classification system for the **Make Data Count (MDC)** competition, demonstrating the application of both machine learning and deep learning concepts to a real-world NLP challenge. The primary objective was to identify whether a dataset mention in a scientific article was **Primary**, **Secondary**, or **Missing**.

### 6.1. Summary of Key Findings

- A robust **data extraction and preprocessing pipeline** was developed using XML parsing and regex-based context identification.
- Exploratory Data Analysis revealed significant **class imbalance**, distinct **text length patterns**, and unique **n-gram distributions** across labels.
- A baseline **TF-IDF + Logistic Regression** model achieved a macro-averaged F1-score of **0.79**, with high accuracy in classifying “Missing” and “Secondary” mentions.
- Visualization of the **confusion matrix** indicated good performance overall, with room for improvement particularly in detecting **Primary** mentions, which were both less frequent and linguistically subtle.

### 6.2. Strengths of Using Deep Learning in Citation Classification

Although this project primarily deployed classical ML techniques, the role of **deep learning in citation classification is transformative**:

- **Transformer-based models** like **BERT** and **RoBERTa** excel at capturing contextual nuances in academic writing—essential for distinguishing between dataset reuse and original data collection.

- **Bidirectional models** (e.g., BiLSTM) can better model the sentence structure and token dependencies.
- Deep learning allows **end-to-end modeling** with minimal manual feature engineering, which improves generalization and scalability.

Future extensions of this work could incorporate **fine-tuned BERT-based architectures** using domain-specific corpora (e.g., SciBERT, BioBERT) for even greater accuracy and semantic understanding.

### 6.3. Limitations and Areas for Improvement

While the current pipeline is effective, several limitations were observed:

- **Limited training examples** for the “Primary” class restrict performance.
- **Ambiguity in short mentions** (e.g., raw DOIs) challenges classifiers, especially without full-text context.
- **Rule-based extraction** can miss complex or nested dataset citations across XML structures.

To address these, future work should explore:

- **Named Entity Recognition (NER)** models to localize dataset names
- Enhanced **context windowing** around mentions
- **Multimodal approaches**, combining citation graphs, metadata, and text features

### 6.4. Generalizability to Other Domains

The methods used here are highly transferable:

- In **medical literature**, similar pipelines could classify references to clinical trials, datasets (e.g., MIMIC-III), or instruments.
- In **legal or patent texts**, the system could tag referenced regulations or prior patents.
- With domain-specific tokenizers and pretrained models, this framework can be **adapted across disciplines**, supporting meta-research, digital libraries, and automated literature analysis at scale.

## 6.5. Final Remarks

This project not only demonstrates the power of AI-driven citation analysis but also reinforces the importance of reproducible, transparent, and scalable NLP solutions for modern scientific communication. Through a blend of structured pipelines, statistical learning, and future-ready architectures, it lays a foundation for more intelligent systems that track, evaluate, and interpret the data-centric evolution of research.

## 7. References

- <sup>3</sup> Devlin, J., Chang, M. W., Lee, K. and Toutanova, K., 2019. *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding*. [online] arXiv. Available at: <https://arxiv.org/abs/1810.04805> [Accessed 19 June 2025].
- <sup>4</sup> Goodfellow, I., Bengio, Y. and Courville, A., 2016. *Deep Learning*. Cambridge, MA: MIT Press. <sup>7</sup>
- Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L. and Stoyanov, V., 2019. *RoBERTa: A Robustly Optimized BERT Pretraining Approach*. [online] arXiv. Available at: <https://arxiv.org/abs/1907.11692> [Accessed 19 June 2025].
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł. and Polosukhin, I., 2017. *Attention is All You Need*. [online] arXiv. Available at: <https://arxiv.org/abs/1706.03762> [Accessed 19 June 2025].
- <sup>12</sup> Beltagy, I., Lo, K. and Cohan, A., 2019. SciBERT: A Pretrained Language Model for Scientific Text. *arXiv preprint* [online] Available at: <https://arxiv.org/abs/1903.10676> [Accessed 19 June 2025].
- Devlin, J., Chang, M. W., Lee, K. and Toutanova, K., 2019. *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding*. *arXiv preprint* [online] Available at: <https://arxiv.org/abs/1810.04805> [Accessed 19 June 2025].
- <sup>9</sup> Habibi, M., Weber, L., Neves, M., Wiegandt, D.L. and Leser, U., 2017. Deep learning with word embeddings improves biomedical named entity recognition. *Bioinformatics*, 33(14), pp.i37–i48.

<sup>5</sup> Lee, J., Yoon, W., Kim, S., Kim, D., Kim, S., So, C.H. and Kang, J., 2020. BioBERT: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, 36(4), pp.1234–1240.

<sup>13</sup> Lo, K., Wang, L.L., Neumann, M., Kinney, R. and Weld, D.S., 2020. S2ORC: The Semantic Scholar Open Research Corpus. *Proceedings of ACL*, pp.4969–4983.

Mayer, R., Kunneman, F., Boon, M. and van den Bosch, A., 2021. Comparing Deep Learning and Machine Learning Models for Scientific Citation Classification. *Journal of Informetrics*, 15(3), p.101188.

Pielke, R., Peiser, L. and Morshed, A., 2022. A hybrid approach for identifying dataset mentions in social science papers. *Data Science Journal*, 21(1), p.23.

<sup>8</sup> Wadden, D., Wennberg, U., Luan, Y. and Hajishirzi, H., 2019. Entity, Relation, and Event Extraction with Contextualized Span Representations. *arXiv preprint* [online] Available at: <https://arxiv.org/abs/1909.03546> [Accessed 19 June 2025].

<sup>19</sup> Yoon, W., Kim, S., Kim, D., Kim, S. and Kang, J., 2019. Chemical-gene relation extraction using recursive neural networks. *Database*, 2019.

<sup>6</sup> Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M. et al., 2011. *Scikit-learn: Machine learning in Python*. *Journal of Machine Learning Research*, 12, pp.2825–2830.

<sup>3</sup> Devlin, J., Chang, M. W., Lee, K. and Toutanova, K., 2019. *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding*. *arXiv preprint*. Available at: <https://arxiv.org/abs/1810.04805> [Accessed 19 June 2025].

Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M. et al., 2011. *Scikit-learn: Machine learning in Python*. *Journal of Machine Learning Research*, 12, pp.2825–2830.

<sup>4</sup> Devlin, J., Chang, M. W., Lee, K. and Toutanova, K., 2019. *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding*. *arXiv preprint*. Available at: <https://arxiv.org/abs/1810.04805> [Accessed 19 June 2025].

Devlin, J., Chang, M. W., Lee, K. and Toutanova, K., 2019. *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding*. *arXiv preprint*. Available at: <https://arxiv.org/abs/1810.04805> [Accessed 19 June 2025].

Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M. et al., 2011. *Scikit-learn: Machine learning in Python*. Journal of Machine Learning Research, 12, pp.2825–2830.

Screenshot of competition

The screenshot shows the Kaggle website interface for the competition 'Make Data Count - Finding Data References'. The left sidebar contains navigation links: Home, Competitions, Datasets, Models, Code, Discussions, Learn, More, Your Work, and a 'views' section with links to 'Make Data Count - FI...', 'Train and Test Datas...', and 'View Active Events'. The main content area displays the 'Leaderboard' tab. At the top of the leaderboard, there is a search bar, a 'Submit Prediction' button, and tabs for Overview, Data, Code, Models, Discussion, Leaderboard, Rules, Team, and Submissions. The leaderboard table lists participants with their rank, name, profile picture, score, number of submissions, and time taken. The participants are sorted by score in descending order. The first participant is 'Alfahesha' with a score of 0.041, 10 submissions, and 20d time. The second is 'Team Lerno' with a score of 0.041, 24 submissions, and 5d time. The third is 'Yash Marathe' with a score of 0.040, 3 submissions, and 1mo time. The fourth is 'krishan abeywardhana' with a score of 0.040, 1 submission, and 1mo time. Below the fourth participant, there is a message: 'Your First Entry! Welcome to the leaderboard!'. The fifth participant is 'Godwin Kofi Klutse' with a score of 0.039, 1 submission, and 1mo time. The sixth is 'YahiaFadi' with a score of 0.039, 10 submissions, and 7d time. The seventh is 'YC' with a score of 0.039, 1 submission, and 12d time.

Rank	Name	Score	Submissions	Time
479	Alfahesha	0.041	10	20d
480	Team Lerno	0.041	24	5d
481	Yash Marathe	0.040	3	1mo
482	krishan abeywardhana	0.040	1	1mo
Your First Entry! Welcome to the leaderboard!				
483	Godwin Kofi Klutse	0.039	1	1mo
484	YahiaFadi	0.039	10	7d
485	YC	0.039	1	12d

ORIGINALITY REPORT

23%	22%	17%	15%
SIMILARITY INDEX	INTERNET SOURCES	PUBLICATIONS	STUDENT PAPERS

PRIMARY SOURCES

1	Submitted to University of Exeter Student Paper	3%
2	Submitted to Institute of Technology, Sligo Student Paper	1%
3	Submitted to University of Derby Student Paper	1%
4	<a href="http://www.frontiersin.org">www.frontiersin.org</a> Internet Source	1%
5	<a href="http://arxiv.org">arxiv.org</a> Internet Source	1%
6	<a href="http://www.bundesbank.de">www.bundesbank.de</a> Internet Source	1%
7	Haneen Deeb, Aldert Vrij, Nicola Palena, Petra Hypšová, Gerges Dib, Sharon Leal, Samantha Mann. "Honesty repeats itself: comparing manual and automated coding on the veracity cues total details and redundancy", Applied Psycholinguistics, 2024 Publication	1%
8	Submitted to Kaplan Professional Student Paper	1%
9	<a href="http://biocreative.bioinformatics.udel.edu">biocreative.bioinformatics.udel.edu</a> Internet Source	1%
10	<a href="http://aigradients.com">aigradients.com</a> Internet Source	1%
11	<a href="http://assets-eu.researchsquare.com">assets-eu.researchsquare.com</a> Internet Source	1%

12	<a href="http://www.biorxiv.org">www.biorxiv.org</a> Internet Source	1 %
13	<a href="http://link.springer.com">link.springer.com</a> Internet Source	1 %
14	<a href="http://www.mdpi.com">www.mdpi.com</a> Internet Source	1 %
15	Submitted to Liverpool John Moores University Student Paper	1 %
16	<a href="http://editorialge.com">editorialge.com</a> Internet Source	1 %
17	Submitted to University College Dublin (UCD) Student Paper	<1 %
18	<a href="http://ijrpr.com">ijrpr.com</a> Internet Source	<1 %
19	<a href="http://hdl.handle.net">hdl.handle.net</a> Internet Source	<1 %
20	<a href="http://aclanthology.org">aclanthology.org</a> Internet Source	<1 %
21	Thangaprakash Sengodan, Sanjay Misra, M Murugappan. "Advances in Electrical and Computer Technologies", CRC Press, 2025 Publication	<1 %
22	Khumukcham Robindro Singh, Nazrul Hoque, Arnab Kumar Maji, Sabyasachi Mondal et al. "Emerging Trends and Future Directions in Artificial Intelligence, Machine Learning, and Internet of Things Innovations - A proceeding of NEIAIS — 2025", CRC Press, 2025 Publication	<1 %
23	<a href="http://opus.lib.uts.edu.au">opus.lib.uts.edu.au</a> Internet Source	<1 %



24	<a href="http://www.coursehero.com">www.coursehero.com</a> Internet Source	<1 %
25	<a href="http://www.public.iastate.edu">www.public.iastate.edu</a> Internet Source	<1 %
26	Submitted to Asia Pacific University College of Technology and Innovation (UCTI) Student Paper	<1 %
27	Submitted to Institute of Research & Postgraduate Studies, Universiti Kuala Lumpur Student Paper	<1 %
28	<a href="http://www.capgemini.com">www.capgemini.com</a> Internet Source	<1 %
29	Submitted to Business School Lausanne Student Paper	<1 %
30	<a href="http://5dok.org">5dok.org</a> Internet Source	<1 %
31	Submitted to Monash University Student Paper	<1 %
32	Submitted to Kingston University Student Paper	<1 %
33	Submitted to UCL Student Paper	<1 %
34	"Chinese Computational Linguistics", Springer Science and Business Media LLC, 2019 Publication	<1 %
35	<a href="http://ecomagazine.com">ecomagazine.com</a> Internet Source	<1 %
36	<a href="http://mental.jmir.org">mental.jmir.org</a> Internet Source	<1 %
37	<a href="http://www.ijcaonline.org">www.ijcaonline.org</a> Internet Source	<1 %

38	<a href="http://www.indusedu.org">www.indusedu.org</a> Internet Source	<1 %
39	<a href="http://bmcmcdinformdecismak.biomedcentral.com">bmcmcdinformdecismak.biomedcentral.com</a> Internet Source	<1 %
40	<a href="http://ai.jmir.org">ai.jmir.org</a> Internet Source	<1 %
41	<a href="http://medwinpublishers.com">medwinpublishers.com</a> Internet Source	<1 %
42	<a href="http://www.ijainn.latticescipub.com">www.ijainn.latticescipub.com</a> Internet Source	<1 %
43	<a href="http://irojournals.com">irojournals.com</a> Internet Source	<1 %
44	<a href="http://medium.com">medium.com</a> Internet Source	<1 %
45	<a href="http://press.amu.edu.pl">press.amu.edu.pl</a> Internet Source	<1 %
46	<a href="http://pure.uva.nl">pure.uva.nl</a> Internet Source	<1 %
47	<a href="http://www.internationalpubls.com">www.internationalpubls.com</a> Internet Source	<1 %

Exclude quotes Off

Exclude matches Off

Exclude bibliography Off