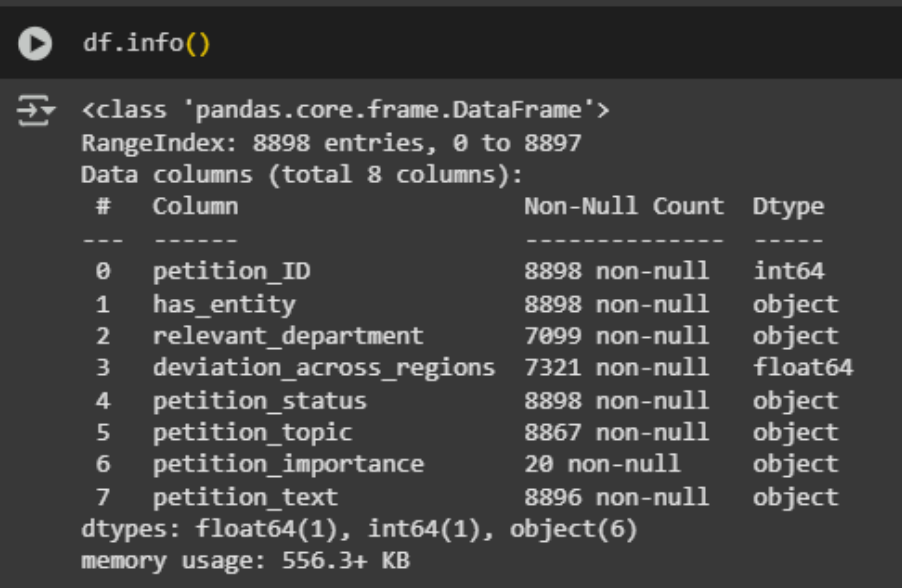**EXECUTIVE SUMMARY**

For Task 1, I implemented a DistilBERT-based model with text preprocessing and entity handling. The model achieved a test accuracy of 99%, exceeding the client's 86% requirement. No overfitting was detected, and all seven classes met the ≤13% misclassification threshold, including the "uk government and devolution" class, which had a 0% misclassification rate.

For Task 2, I developed a semi-supervised DistilBERT model and a RandomForest model for classifying petition importance. Compared to a test accuracy of 17.9% for the majority class baseline, the DistilBERT model achieved a test accuracy of 89.74%, while the RandomForest classifier achieved a test accuracy of 97.43%. The ethical analysis found risks associated with transparency, representation bias, and deprioritizing critical petitions and suggested ways to mitigate them. The models showed no overfitting and demonstrated strong performance across both the "important" and "not_important" categories.

## 1. DATA EXPLORATION AND ASSESSMENT

The "pandas library", a powerful tool for loading CSV files, handling data, and manipulating data (Gupta and Bagchi, 2024) was used to load and transform the dataset into a dataframe within the Google Colab environment. When the dataframe was examined using "info()" and "isnull().sum()", it was found to have 8898 rows and 8 columns with both numerical and categorical data, and there were missing values (see Figure 1.1). Missing data in machine learning negatively impacts model accuracy, introduces bias, reduces sample size, and can distort feature relationships, ultimately leading to unreliable and potentially misleading results (Emmanuel *et al.*, 2021). The 'petition_importance' column showed an initial label distribution of 'unknown' 88.7%—8,878 petitions, 'important' 0.1%—11 petitions, and 'not_important' 0.1%—9 petitions. This imbalance necessitated careful labeling strategies and semi-supervised learning approaches (Task 2).

```
df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 8898 entries, 0 to 8897
Data columns (total 8 columns):
 #   Column                 Non-Null Count  Dtype
---  ------                 --------------  -----
 0   petition_ID            8898 non-null   int64
 1   has_entity             8898 non-null   object
 2   relevant_department    7099 non-null   object
 3   deviation_across_regions  7321 non-null  float64
 4   petition_status        8898 non-null   object
 5   petition_topic         8867 non-null   object
 6   petition_importance    20 non-null     object
 7   petition_text          8896 non-null   object
dtypes: float64(1), int64(1), object(6)
memory usage: 556.3+ KB
```

**Figure 1.1: Summary Information of the Petition Dataset**

When the 'petition_topic' column was examined using "value_counts(normalize=True)," duplicate values were found, increasing the number of topics to 14 (see Figure 1.2). John (2024) reports that duplicate values in a dataset used for machine learning can lead to biased models, overfitting, and inflated accuracy, making it crucial to identify and address them. The topics showed some class imbalance, particularly with the "london" category being significantly underrepresented (see Figure 1.2). Tolstaya *et al.* (2025) report the challenges of class imbalance in a machine learning model, one of which is model bias.
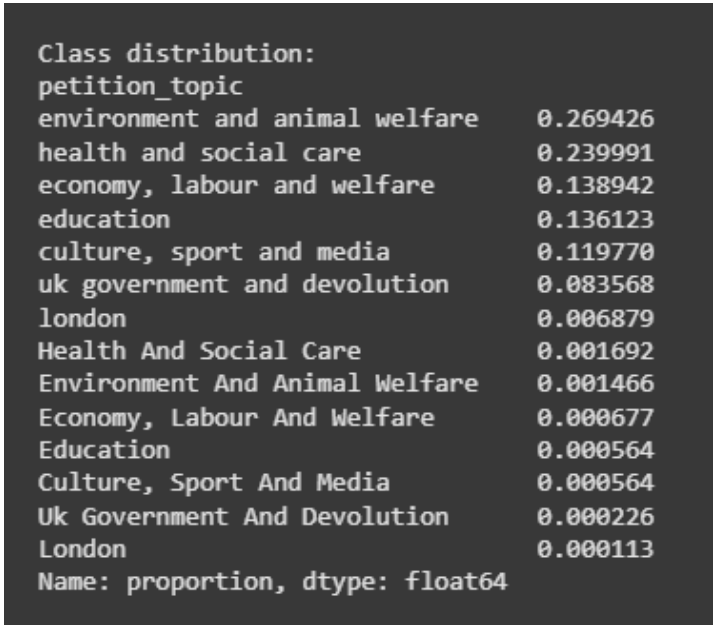
```
Class distribution:
petition_topic
environment and animal welfare    0.269426
health and social care            0.239991
economy, labour and welfare       0.138942
education                         0.136123
culture, sport and media          0.119770
uk government and devolution      0.083568
london                            0.006879
Health And Social Care            0.001692
Environment And Animal Welfare    0.001466
Economy, Labour And Welfare       0.000677
Education                         0.000564
Culture, Sport And Media          0.000564
Uk Government And Devolution      0.000226
London                            0.000113
Name: proportion, dtype: float64
```

**Figure 1.2: Class Distribution of Petition Topic**

The length of the petition text was also examined, and the results showed variation across petitions, with some very short and some very long texts (see Figure 1.3). Feldman (2018) reports on the importance of a petition's length. This insight informed decisions about text preprocessing and manual labeling implemented in Task 2.
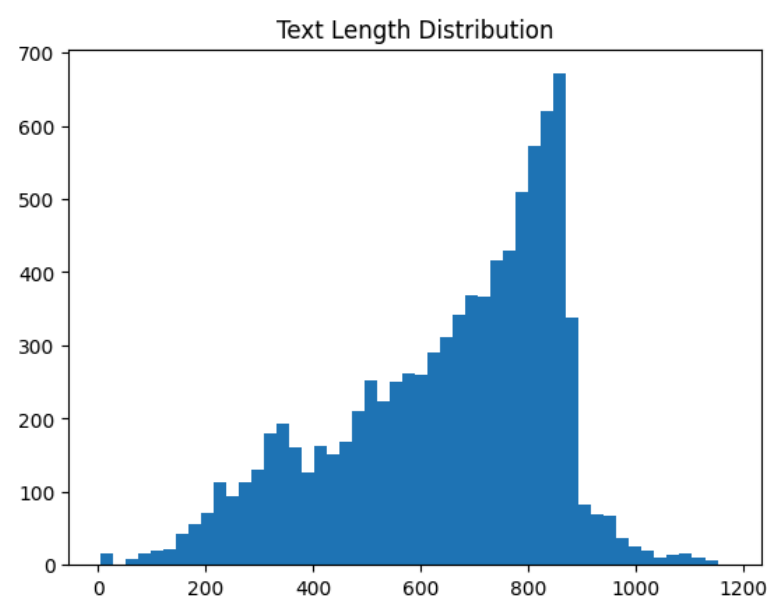


**Figure 1.3: Visualisation of Petition Text Distribution**

## 2.  DATA SPLITTING AND CLEANING

For Task 1, the dataset was prepared as follows:

**i. Handling Missing Values:** This was necessary because missing target values would prevent supervised learning algorithms from functioning optimally (Tsai and Hu, 2022). This approach minimizes potential biases and inaccuracies that imputation might introduce. As noted in a discussion on best practices for handling missing data, dropping rows with missing values is advisable when the number of such rows is minimal compared to the dataset's size (LinkedIn, 2024).

**ii. Value Harmonisation:** The values in the 'petition_topic' column were harmonised by converting them to lowercase using 'str.lower()' to prevent duplicates that arise due to case sensitivity (Upadhye, 2020; Ardalan *et al.*, 2022).

**iii. Data Splitting:** The data was split using StratifiedKFold with 5 folds to ensure a representative distribution of classes across training and validation sets.

**iv. Class Imbalance Handling:** Class weights were computed using 'compute_class_weight' to address the imbalance in the 'petition_topic' distribution. Ng (2025) reports that implementing 'compute_class_weight' ensures that the model does not become biased toward the majority class, thereby improving its generalization across all classes.

For Task 2, the following steps were taken:

**i. Handling Missing Values:** Missing values in feature variables ('relevant_department', 'deviation_across_regions', 'petition_status', 'petition_topic', 'petition_importance') were filled with the string 'unknown'. Perez-Lebel *et al.* (2022) report that treating missing values as 'unknown' categories can be beneficial, especially when the lack of data carries meaningful information, thus maintaining the dataset's integrity and ensuring that the model is aware of missingness, which can be informative, improving model performance.

**ii. Data Splitting:** The data was divided into labeled and unlabeled portions. The labeled data was then split into training and test sets using train_test_split with stratification to maintain class distribution (Galarnyk and Whitfield, 2025).

**iii. Semi-Supervised Learning Setup:** The unlabeled data was used with the labeled data to train the model, allowing it to leverage patterns from the larger dataset while still having ground truth for supervision (Dong *et al.*, 2024). The technique of manually labeling unlabeled data is known as pseudo-labeling, which leverages patterns from the unlabeled data, enhancing the model's performance (Dong *et al.*, 2024).

**Updated dataset characteristics:** 8867 observations (instead of the original 8898), no null values for the 'petition_topic' (instead of the original 31 null values), the values in 'petition_topic' reduced to 7 from 14.

## 3.  DATA ENCODING

For Task 1, the following encoding steps were performed:

**i. Text Encoding:** The 'petition_text' feature underwent the following extensive preprocessing implementation:

- Special characters and numbers were removed using regular expressions. Ganesan (2019) reports that this process is known as noise removal, and it is essential for retaining only the textual content relevant to the task.

- Each text observation was converted to lowercase.
- Tokenization was performed using NLTK's word_tokenize. This fundamental step in NLP helps facilitate the analysis of individual text components (Deepanshi, 2025).
- Stopwords were removed using NLTK's English stopwords list. This step enhances the model's ability to learn from the data (Agrawal, 2021).
- Lemmatization was applied using NLTK's WordNetLemmatizer. Lemmatization ensures that different forms of a word are treated as a single entity, improving the consistency of the data (Agrawal, 2021).

The cleaned text was then encoded using DistilBERT's tokenizer, which converts text into token IDs and other inputs required by the DistilBERT model (Sanh *et al.*, 2019).

**ii. Feature Engineering and Value Encoding:** The 'has_entity' feature was split into separate columns, creating three binary features that could be directly used in the model. Milwaukee School of Engineering (2024) reports on several advantages of feature engineering to machine learning models, while Matteucci *et al.* (2023) highlight the significance of encoding categorical data for machine learning.

**iii. Model-Specific Encoding:** Sanh *et al.* (2019) report that the DistilBERT model processes text through its embedding layer, which converts token IDs into contextual embeddings. Sanh *et al.* (2019) highlight that for the model to comprehend and process natural language, the embedding layer is essential to this transformation.

For Task 2, the encoding approach was slightly different to accommodate the semi-supervised learning setup:

**i. Text Encoding:** The text preprocessing procedure of Task 1 was applied, followed by TF-IDF vectorization for the RandomForest baseline model. Das *et al.* (2023) revealed that the Random Forest classifier achieved high accuracy when utilizing TF-IDF features. For the DistilBERT model, the same DistilBERT tokenizer was used as in Task 1.

**ii. Value Encoding:** Unlike Task 1, the 'has_entity' feature remained a single column; however, a lambda function was applied to convert the associated values from 'YES/NO' to 1/0. There is no justification for this variation other than experimentation.

**ii. Numeric Feature Standardization:** Numeric features were standardized using Scikit-Learn's StandardScaler. Sánchez (2024) reports that the StandardScaler standardises features by removing the mean and scaling to unit variance, ensuring that each feature contributes equally to the model's performance.

**iii. Categorical Feature Encoding:** Categorical features were one-hot encoded using Scikit-Learn's OneHotEncoder. This transformation allows the model to consider each category as a separate dimension (Matteucci *et al.*, 2023).

**iv. Feature Combination:** For the RandomForest model, the ColumnTransformer (embedded with StandardScaler and OneHotEncoder) was used to transform numeric and categorical features. Rabiller (2024) highlights the effectiveness of integrating TF-IDF text features with numeric and categorical data to improve model performance. For the DistilBERT model, numeric features were normalized and concatenated with the text embeddings before being passed to the classifier layer.

## 4.    TASK 1: TOPIC CLASSIFICATION

### a.    Model Building

The model operates as follows:

- Textual input is processed through DistilBERT's transformer layers to generate contextualized embeddings (Sanh *et al.*, 2019).
- Entity features are passed through a separate fully connected (dense) layer to extract meaningful representations.
- The outputs from both components are concatenated before classification.
- A final, fully connected layer maps the combined representation to the output classes.

This architecture enables the model to leverage unstructured and structured information for improved classification accuracy.

Figure 4.1 presents the hyperparameters used:

| Hyperparameter | Value | Description |
| --- | --- | --- |
| MAX_LEN | 256 | Maximum sequence length for text input |
| BATCH_SIZE | 16 | Number of samples per training batch |
| EPOCHS | 3 | Total number of training iterations over data |
| Learning Rate | 2e-5 | Initial learning rate for **AdamW** optimizer |
| Dropout Rate | 0.1 | Dropout probability for regularization |
| Entity Layer Units | 32 | Number of neurons in entity processing layer |
| Classifier Layer Input Size | 768 + 32 | Combined size of text and entity features |
| Classifier Layer Output Size | 7 | Number of target classes |

**Figure 4.1: Hyperparameters used for the DistilBERT model**

The model was trained using the AdamW optimizer with weight decay to enhance convergence and reduce overfitting (Loshchilov and Hutter, 2019). Cross-entropy loss combined with computed class weights addressed the dataset's imbalance (Cui *et al.*, 2019), while a dropout of 0.1 and gradient clipping ensured training stability (Wen *et al.*, 2025). A batch size of 16 and a learning rate of 2e-5 maintained efficient memory use and controlled convergence (Devlin *et al.*, 2019). Employing 5-fold cross-validation further validated the model's robustness across varied data splits (Madiha *et al.*, 2025), ultimately producing a well-generalized classifier that effectively integrates textual and structured features.

The selection of hyperparameters was informed by the guidance provided in Sanh *et al.* (2019) on DistilBERT. Based on these findings, I adopted key settings—such as a maximum sequence length of 256, a learning rate of 2e-5, and a batch size of 16. The model was initially evaluated on a held-out validation set using overall accuracy as the primary metric

(see Figure 4.2), and a 5-fold cross-validation approach was subsequently employed to ensure robustness across different data splits (see Figure 4.2).
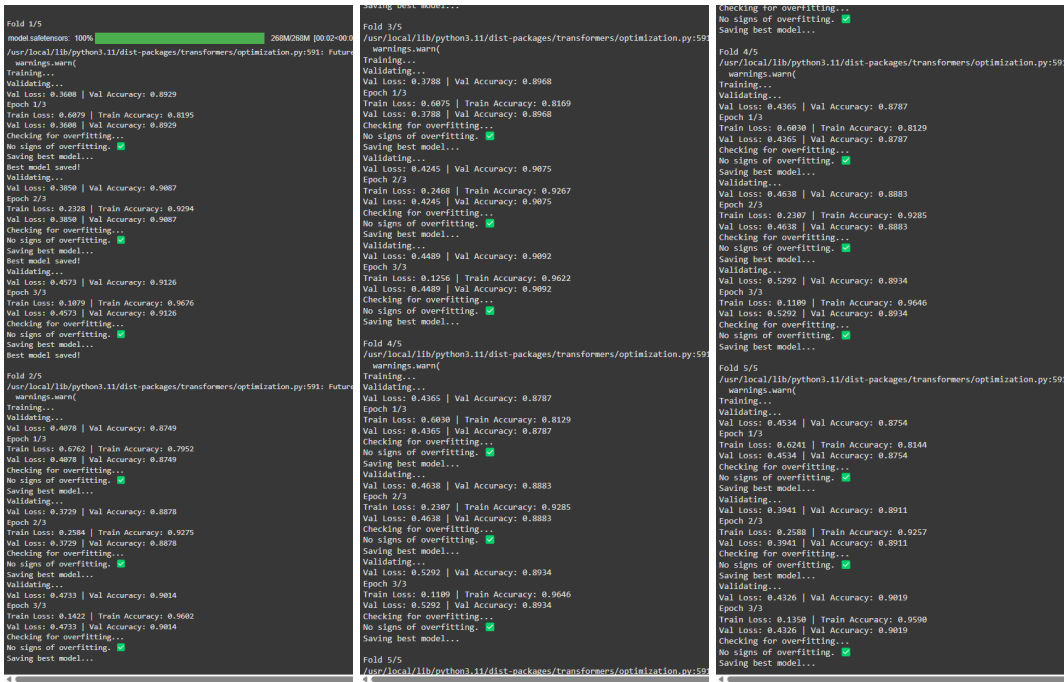


**Figure 4.2: DistilBERT 5-Fold Cross-Validation Summary**

**Justification**
- **Efficient Language Understanding:** DistilBERT was selected for its proven efficiency and strong language comprehension, making it well-suited for processing short petition texts without the heavy computational cost of larger models.
- **Non-empirical research:** A review of studies on classification revealed DistilBERT's ability to capture semantic relationships in text (Sanh *et al.*, 2019), making it a suitable choice for this project.
- **Empirical Validation:** Experimental results, including cross-validation and performance metrics, confirmed that these design choices significantly enhanced accuracy and generalization.

b. **Model Evaluation**

The model's performance was evaluated using multiple metrics to provide a comprehensive understanding of its capabilities (see Figure 4.3):
- Accuracy: 99% on the test set
- Precision: Ranged from 0.92 to 1.00 across classes
- Recall: Ranged from 0.96 to 1.00 across classes
- F1-Score: Ranged from 0.96 to 1.00 across classes

**Figure 4.3: Classification Report**

The confusion matrix revealed strong performance (see Figures 4.3 and 4.4). Notably:
- The model achieved 99% test accuracy, exceeding the client's requirement
- The model showed no signs of overfitting
- All 7 classes met the client misclassification threshold of ≤13%
- The "uk government and devolution" class showed 0% misclassification
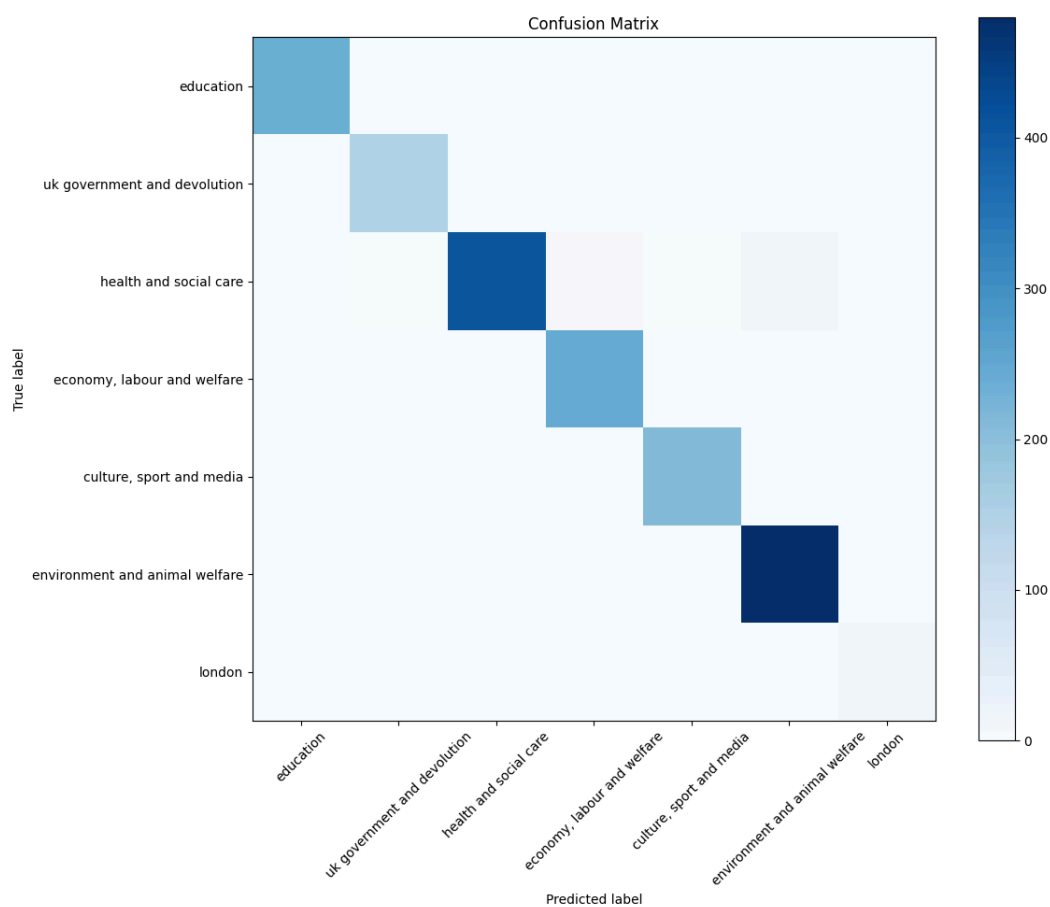


**Figure 4.4: Confusion Matrix**

The classification report showed (see Figure 4.3):
- Macro average precision, recall, and F1-score of 0.98, 0.99, and 0.98, respectively
- Weighted average precision, recall, and F1-score of 0.99 across all classes

Overall, the model demonstrated exceptional performance across all evaluation metrics, significantly exceeding the client's success criteria.

### c. Conclusion

All four of the client's success conditions are met by the model:
- Test accuracy of 99% exceeds the required 86%
- No overfitting detected; validation accuracy remained within 1% of training accuracy
- All seven classes meet the ≤13% misclassification threshold, with misclassification rates ranging from 0% to 11%
- The "uk government and devolution" class achieved 0% misclassification, well below the 9% requirement

I recommend that the client use the macro F1-score to evaluate the model performance. Since the dataset is imbalanced, this metric ensures that all classes are weighted equally, providing a better measure of overall model effectiveness.

## 5. TASK 2: PETITION IMPORTANCE CLASSIFICATION PROTOTYPE

### a. Ethical discussion

The task of predicting petition importance in the UK context presents significant ethical challenges that require careful consideration. Using the Data Hazards Labels framework, we identify three primary risks: bias amplification, automation surprise, and incompleteness. Historical biases in petition engagement could lead the model to deprioritize petitions from marginalized communities or on contentious topics, potentially reinforcing systemic inequalities. The risk of automation surprise arises from potential over-reliance on algorithmic decisions, which might overlook critical petitions without human oversight, thereby affecting democratic engagement. Incompleteness results from potential representation gaps in the training data, particularly for petitions from rural areas or less mainstream issues. These risks are compounded by the subjective nature of determining petition importance, which could lead to framing traps where complex societal issues. To mitigate these concerns, implement bias audits, maintain human oversight for contentious cases, and continuously monitor performance across demographic and topic-based subsets to prevent harm to democratic processes in the UK.

### b. Data labelling

**Labelling Criteria:** I identified features with statistically significant correlations among the labelled 'petition_importance' observations. Features with correlation coefficients above 0.3 were considered strong indicators of importance, while those below 0.3 were considered weak indicators. Then, I used keywords (in the 'petition_text' column) like "urgent," "critical," and "emergency" to identify important petitions.

**Labeling Process:** Petitions containing keywords or from high-correlation departments/topics were labeled "important." While petitions with short text (<100 words) and no high-correlation features were labeled "not_important." See Figures 5.1 and 5.2 for statistics and labelling example texts, respectively.



| petition_importance | count |
|---|---|
| unknown | 8509 |
| important | 317 |
| not_important | 72 |

**Figure 5.1: Final Label Statistics**



```
df[df['petition_importance'] == 'important']['petition_text'][1]
```

'Lock NHS pay rises to increases in pay for MPs, backdated to 2015. Whenever the salaries of MPs rise, the salaries of every NHS staff member should rise by the same proportion. This should be backdated to 2015 to recognise the sacrifices of NHS staff over the pandemic and recognise the large real terms pay cut NHS staff have faced over many years. Between April 2015 and April 2020 MPs salaries have increased by £14,872. Over the same time period, most NHS staff have had real terms pay cuts.\nThis is because when IPSA decides MPs' pay should rise, it happens. Whereas when the independent organisations that determine NHS staff pay release their recommendations, the Government can ignore it.\n \nThis is emblematic of a Government that believes there's one rule for them and another for everyone else.'

```
df[df['petition_importance'] == 'not_important']['petition_text'][4]
```

'Make bank holidays a holiday. Make bank holidays a holiday please because i want bank holidays to be a holiday'

**Figure 5.2: Labeling Examples ('important' and 'not_important')**

I have a moderate level of confidence in my scientific labeling strategy because of the
- Subjective nature of importance determination
- Potential for correlation without causation
- Limited initial labeled examples for correlation analysis

c. **Model building and evaluation**

The model utilized the following combination of text and numeric features:
```
'petition_text', 'deviation_across_regions', 'petition_status',
'relevant_department', 'petition_topic', 'petition_text_length',
'has_entity'
```

These features were selected based on their correlation with `petition_importance`, identified through Pearson correlation analysis on the initial labeled dataset.

The final model used a custom neural network combining DistilBERT, a linear layer, a combination of both representations, and a classifier layer. Below are the hyperparameters used:

| Hyperparameter | Value | Description |
|---|---|---|
| MAX_LEN | 128 | Maximum sequence length for text input |
| BATCH_SIZE | 32 | Batch size for training |
| EPOCHS | 5 | Number of training epochs |
| Learning Rate | 3e-5 | Initial learning rate for AdamW optimizer |
| Dropout Rate | 0.3 | Dropout rate for regularization |
| Numeric Layer Units | 32 | Number of units in numeric processing layer |
| Classifier Layer Input Size | 768 + 32 | Combined size of text and numeric features |
| Classifier Layer Output Size | 2 | Number of classes |

**Figure 5.3: Hyperparameters**

The advanced techniques used include:
- Semi-supervised Learning
- Dynamic Learning Rate Adjustment
- Gradient Clipping
- Cross-Entropy Loss combined with class weights to handle imbalance
- Stratified k-Fold Cross-Validation

**Justifications**
- DistilBERT was used for its balance between computational efficiency and strong text understanding capabilities.
- Semi-supervised learning addressed the limited labeled data while leveraging the full dataset for better generalization.
- Regularization Techniques prevented overfitting despite the complex architecture and limited labeled data.

The model showed strong performance when compared to the baseline (see the results below):
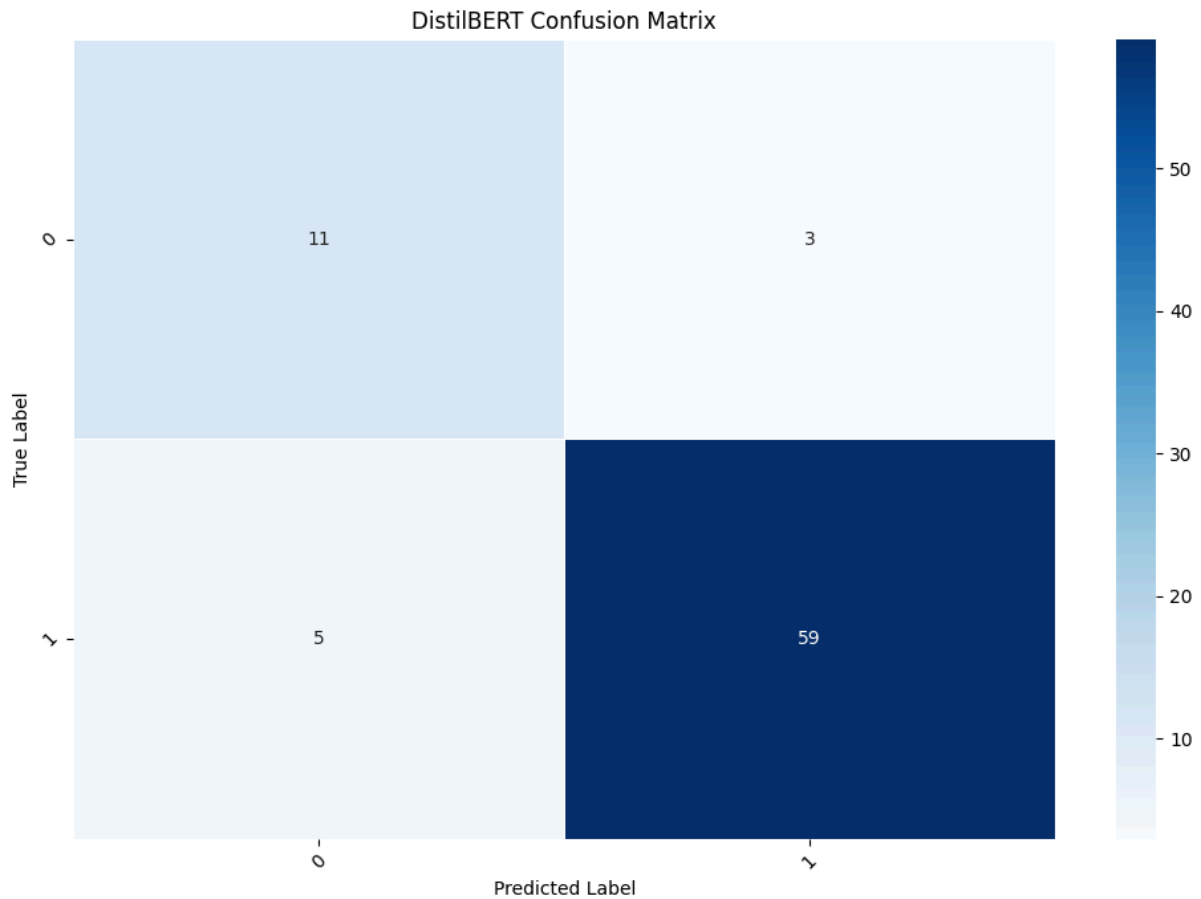
**Figure 5.4: Final model Confusion Matrix**



```
Classification Report for DistilBERT Model:
              precision    recall  f1-score   support

           0       0.69      0.79      0.73        14
           1       0.95      0.92      0.94        64

    accuracy                           0.90        78
   macro avg       0.82      0.85      0.83        78
weighted avg       0.90      0.90      0.90        78
```

**Figure 5.5: Final model Classification Report**

| Model | Accuracy | F1-score |
| --- | --- | --- |
| Baseline (Majority Class) | 0.179 | 0.739 |
| DistilBERT Model | 0.897 | 0.900 |

**Figure 5.6: Comparing Model Performances**

The model meets the requirement of 86% test accuracy and avoids being deemed unethical even with the limited labelled data.

**d. Task 2 Conclusions**

- The model is successful according to the client's definition of success. It performs significantly better than the majority class baseline (89.7% vs. 17.9% accuracy). While ethical considerations are significant, mitigation strategies can be properly applied.
- The task is well-framed given the available data, although the limited number of labeled examples for "not_important" petitions presents some challenges. The application of semi-supervised learning appropriately addresses the challenge of labeled data.
- The top suggestion for improvement is to obtain more labeled data, particularly for the "not_important" class, to address class imbalance and enhance the model's performance on that category.

**6. SELF-REFLECTION**

The most challenging aspect was balancing model performance with ethical considerations while managing limited labeled data. If starting over, I would allocate more time to initial data exploration and spend more time on ethical risk mitigation.

**REFERENCES**

Agrawal, R. (2021) 'Must known techniques for text preprocessing in NLP', Analytics
Vidhya, 15 October. Available at:
https://www.analyticsvidhya.com/blog/2021/06/must-known-techniques-for-text-prepro
cessing-in-nlp/ (Accessed: 16 March 2025).

Ardalan, A., Paulsen, D., Saini, A.S., Cai, W. and Doan, A. (2022) 'Toward data cleaning
with a target accuracy: A case study for value normalization', in 2022 IEEE
International Conference on Big Data (Big Data), December, pp. 3975–3981. IEEE.
Available at: https://doi.org/10.48550/arXiv.2101.05308

Cui, Y., Jia, M., Lin, T.Y., Song, Y. and Belongie, S. (2019) 'Class-balanced loss based on
effective number of samples', in Proceedings of the IEEE/CVF Conference on
Computer Vision and Pattern Recognition, pp. 9268–9277. Available at:
https://doi.org/10.1109/CVPR.2019.00949

Das, M. and Alphonse, P.J.A. (2023) 'A comparative study on tf-idf feature weighting
method and its analysis using unstructured dataset', arXiv preprint arXiv:2308.04037.
Available at: https://doi.org/10.48550/arXiv.2308.04037

Deepanshi (2025) 'Text preprocessing in NLP with Python codes', Analytics Vidhya, 10
March. Available at:
https://www.analyticsvidhya.com/blog/2021/06/text-preprocessing-in-nlp-with-python-
codes/ (Accessed: 16 March 2025).

Devlin, J., Chang, M.-W., Lee, K. and Toutanova, K. (2019) 'BERT: Pre-training of deep
bidirectional transformers for language understanding', in Proceedings of the 2019
Conference of the North American Chapter of the Association for Computational
Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers).
Minneapolis, Minnesota: Association for Computational Linguistics, pp. 4171–4186.
Available at: https://doi.org/10.18653/v1/N19-1423

Dong, X., Ouyang, T., Liao, S., Du, B. and Shao, L. (2024) 'Pseudo-labeling based practical
semi-supervised meta-training for few-shot learning', IEEE Transactions on Image
Processing. Available at: https://doi.org/10.48550/arXiv.2207.06817

Emmanuel, T., Maupong, T., Mpoeleng, D., Semong, T., Mphago, B. and Tabona, O. (2021)
'A survey on missing data in machine learning', Journal of Big Data, 8, pp. 1–37.
Available at: https://doi.org/10.1186/s40537-021-00516-9

Feldman, A. (2018) 'Empirical SCOTUS: Follow the experts in framing petitions for cert',
SCOTUSblog, 19 November. Available at:
https://www.scotusblog.com/2018/11/empirical-scotus-follow-the-experts-in-framing-p
etitions-for-cert/ (Accessed: 15 March 2025).

Galarnyk, M. and Whitfield, B. (2025) 'Train test split: what it means and how to use it',
Built In, 3 February. Available at: https://builtin.com/data-science/train-test-split
(Accessed: 15 March 2025).

Ganesan, K. (2019) 'All you need to know about text preprocessing for NLP and Machine
Learning', KDnuggets, 5 April. Available at:
https://www.kdnuggets.com/2019/04/text-preprocessing-nlp-machine-learning.html
(Accessed: 16 March 2025).

Gupta, P. and Bagchi, A. (2024) 'Introduction to Pandas', in Essentials of Python for
Artificial Intelligence and Machine Learning. Synthesis Lectures on Engineering,
Science, and Technology. Cham: Springer. Available at:
https://doi.org/10.1007/978-3-031-43725-0_5

John, A.C. (2025) 'Mastering duplicate data management in machine learning for optimal
model performance', DagsHub Blog, 15 January. Available at:
https://dagshub.com/blog/mastering-duplicate-data-management-in-machine-learning-f
or-optimal-model-performance/ (Accessed: 15 March 2025).

LinkedIn (2024) 'Best practices for dealing with missing values and imputation'. Available
at:
https://www.linkedin.com/advice/0/what-some-best-practices-dealing-missing-values-i
mputation (Accessed: 15 March 2025).

Loshchilov, I. and Hutter, F. (2019) 'Decoupled weight decay regularization', in 7th
International Conference on Learning Representations (ICLR). Available at:
https://dblp.org/rec/conf/iclr/LoshchilovH19.html

Madiha, M., Bukhari, I.E.R., Khan, G.S.A. and Kumar, S. (2025) 'CCR-ML: A Machine Learning Approach for Credit Card Risk Classifications', Journal of Computer Sciences and Informatics, 2(1), pp. 57–57. Available at: http://dx.doi.org/10.5455/JCSI.20241231072950

Matteucci, F., Arzamasov, V. and Böhm, K. (2023) 'A benchmark of categorical encoders for binary classification', Advances in Neural Information Processing Systems, 36, pp. 54855–54875. Available at: https://doi.org/10.48550/arXiv.2307.09191

Milwaukee School of Engineering (2024) 'The importance of feature engineering in machine learning', Milwaukee School of Engineering, 3 April. Available at: https://online.msoe.edu/engineering/blog/importance-of-feature-engineering-in-machine-learning (Accessed: 16 March 2025).

Ng, C. (2025) 'How to fine-tune DistilBERT for emotion classification', Towards Data Science, 18 February. Available at: https://towardsdatascience.com/how-to-fine-tune-distilbert-for-emotion-classification/ (Accessed: 15 March 2025).

Perez-Lebel, A., Varoquaux, G., Le Morvan, M., Josse, J. and Poline, J.B. (2022) 'Benchmarking missing-values approaches for predictive models on health databases', GigaScience, 11, p. giac013. Available at: https://doi.org/10.1093/gigascience/giac013

Sanh, V., Debut, L., Chaumond, J. and Wolf, T. (2019) 'DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter', arXiv preprint arXiv:1910.01108. Available at: https://doi.org/10.48550/arXiv.1910.01108

Sánchez, M.S. (2024) 'Optimization of diamond price prediction strategies using machine learning techniques', Journal of Economics and International Finance, 16(3), pp. 28–38. Available at: https://academicjournals.org/journal/JEIF/article-full-text-pdf/3F763DF72446.pdf (Accessed: 16 March 2025).

scikit-learn developers (no date) 'sklearn.model_selection.StratifiedKFold', scikit-learn 1.6.1 documentation. Available at: https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.StratifiedKFold.html (Accessed: 15 March 2025).

Tolstaya, E., Tichy, A., Paris, S. and Schwendicke, F. (2025) 'Improving machine learning-based bitewing segmentation with synthetic data', Journal of Dentistry, p. 105679. Available at: https://doi.org/10.1016/j.jdent.2025.105679

Tsai, C.F. and Hu, Y.H. (2022) 'Empirical comparison of supervised learning techniques for missing value imputation', Knowledge and Information Systems, 64(4), pp. 1047–1075. Available at: https://doi.org/10.1007/s10115-022-01661-0

Upadhye, A. (2020) 'A comprehensive survey of text data cleaning techniques: challenges, methods, and best practices', Journal of Scientific and Engineering Research, 7(8), pp. 205–210. Available at: https://jsaer.com/download/vol-7-iss-8-2020/JSAER2020-7-8-205-210.pdf

Wen, X., Zhao, B., Elezi, I., Deng, J. and Qi, X. (2025) '"Principal Components" Enable A New Language of Images', arXiv preprint arXiv:2503.08685. Available at: https://doi.org/10.48550/arXiv.2503.08685