

Deep Learning for Spam Email Classification: Optimizing Pre-Trained Models with Compression Techniques.

Git Repository : https://github.com/snirdav/project_group_9/tree/main

Abstract

This article explores advanced techniques for spam email classification using pre-trained models, XLNet and RoBERTa, from the Hugging Face library. Various model compression techniques, including knowledge distillation, quantization, and pruning, were employed to enhance efficiency. Despite the lack of significant improvement from compression methods, XLNet outperformed RoBERTa, establishing it as the preferred model for further optimization.

1. Introduction

Spam email classification is a critical task in natural language processing (NLP) aimed at accurately identifying and filtering out unwanted emails. Traditional models have achieved significant success, but with the increasing size and complexity of data, there is a growing need for more efficient alternatives. This study focuses on compressing models using knowledge distillation, pruning, and quantization, alongside evaluating the effectiveness of pre-trained transformer models like XLNet and RoBERTa. The effectiveness of XLNet in this study is consistent with the findings of Yang et al. (2019), who introduced XLNet as a powerful transformer model that captures bidirectional context using permutation-based language modeling.

2. Methodology

2.1 Data Processing

The dataset used for this project consists of labeled email data with features derived from the email content. Data preprocessing involved cleaning the text, tokenizing, and transforming it into a format suitable for model training.

2.2 Class Structure and Implementation

The project is structured around several key classes, each handling specific aspects of model training and evaluation:

TrainingPipelineManager: This class orchestrates the entire process, managing data ingestion, model initialization, training, and evaluation. It also coordinates post-training operations, such as model compression and logging metrics. The recent enhancements include expanded capabilities for managing pruning, quantization, and knowledge distillation methods, which are vital for optimizing model size and efficiency.

ModelTrainer: Responsible for handling the core training routines, this class manages model setup, evaluation, and performance logging. It now includes advanced methods for pruning (`prune_model()`), quantization (`quantize_model()`), and training a student model (`train_student()`), enabling it to effectively reduce model size while maintaining accuracy.

TransformerModelWrapper: This class encapsulates transformer models like XLNet and RoBERTa, managing forward passes, training steps, and validation within the PyTorch Lightning framework. It ensures that these models are trained efficiently and evaluated thoroughly.

DataPreprocessing: Expanded to handle complex data preprocessing tasks, this class includes methods like `split_data()`, `tokenize_data()`, and `prepare_data()`. These enhancements ensure that data is efficiently prepared, tokenized, and ready for model training, reducing potential bottlenecks in the training pipeline.

ImprovedXGBoostClassifier: Employed towards the end of the pipeline, this class handles metadata-specific tasks. However, it showed lower performance compared to transformer models and was not included in the final ensemble.

- the models were evaluated on both validation and test datasets. Below are the chosen model's results:

XLNet Model (Default Configuration):

- Validation Loss: 0.022162 (Val) / 0.042439 (Test)
- ROC AUC: 0.993593 (Val) / 0.988681 (Test)
- F1-Score: 0.994319 (Val) / 0.990117 (Test)
- Accuracy: 0.994307 (Val) / 0.990128 (Test)

- **ImprovedXGBoostClassifier:** XGBoost was employed at the end, specifically designed for handling metadata. The XGBoost model showed an accuracy of 86%, which was not sufficient to be combined into the XSLNet ensemble model.

2.3 Pipeline Overview

The pipeline begins with the data processing stage, where raw email data is cleaned and tokenized. The processed data is then passed to the TrainingPipelineManager, which initializes the model (XSLNet, RoBERTa, or XGBoost) and coordinates the training and evaluation process.

This study demonstrates the effectiveness of XLNet in spam email classification, outperforming RoBERTa in all evaluated metrics. The use of checkpoints and PKL files was crucial in saving time and preventing loss of progress when the system crashed, especially during the complex training and evaluation cycles. XLNet's lower validation loss and higher accuracy and F1-scores establish it as the superior model for this application. Traditional models like XGBoost, Random Forests, and Neural Networks, although trained only on metadata and not on the email body text, were still outperformed by the transformer-based models in this specific context

1. **Data Ingestion:** Load and preprocess the dataset.
2. **Model Initialization:** Set up the model architecture and parameters.
3. **Training Loop:** Train the model using the training dataset while monitoring performance on the validation set.
4. **Evaluation:** Evaluate the model on the test set and log performance metrics.
5. **Model Compression:** Apply knowledge distillation and quantization (if applicable).
6. **Final Evaluation:** Assess the compressed model's performance on the test set.

2.4 Important Methods Used

Fraud detection is a critical application in various domains, including financial services and cybersecurity. In spam email classification, detecting fraudulent or malicious emails involves identifying patterns and anomalies that distinguish spam from legitimate communication. Techniques such as anomaly detection, clustering, and machine learning models are employed to enhance detection accuracy.

For this research we tried 2 pre-trained models from huginFace library. XSLNet and RoBERTa are both advanced transformer-based models designed for NLP tasks.

- **XSLNet:** Introduced by Yang et al., XSLNet is a generalized autoregressive pretraining method for language understanding. It overcomes the limitations of BERT by capturing bidirectional contexts using permutation-based language modeling. XSLNet is known for its robustness and effectiveness in

various NLP benchmarks, making it suitable for complex tasks like spam classification.

- **RoBERTa:** developed by Facebook AI, is an optimized version of BERT with modifications in training data size, batch size, and learning rates. It is designed to handle large datasets efficiently and has shown strong performance in multiple NLP tasks. However, in our experiments, XSLNet outperformed RoBERTa, making it the preferred model.

Implementation Details: Our implementation leveraged the HuggingFace library, a popular framework for working with transformer models. The process involved:

1. **Tokenization:** Converting raw text into tokens suitable for model input using the XSLNetTokenizer. For longer texts, we split and truncated some entries to ensure they fit within the tokenizer's limits."
2. **Model Setup:** Initializing the XSLNet model via XSLNetModel.
3. **Training:** The ModelTrainer class handled the training process, including forward passes, loss calculation, and backpropagation.
4. **Evaluation:** The model's performance was assessed using a few standard metrics for classifications tasks.

2.5 Evaluation Methods

To evaluate the performance of our models, we used several key metrics:

- **F1-Score:** A harmonic mean of precision and recall, providing a balanced measure of a model's accuracy, particularly in imbalanced datasets.
- **ROC/AUC:** The Receiver Operating Characteristic (ROC) curve plots the true positive rate against the false positive rate, while the Area Under the Curve (AUC) measures the model's ability to discriminate between classes.

2.6 Compressions Methods

In an effort to reduce the model size without sacrificing performance, we applied the following techniques:

1. **Knowledge Distillation:** Transferring knowledge from a large "teacher" model (XSLNet) to a smaller "student" model (NN), aiming to retain most of the teacher's performance. By getting the output of the pre-trained teacher and adjusting to it?! This method follows the approach outlined by Hinton, Vinyals, and Dean (2015) in their foundational work on knowledge distillation, where knowledge is transferred from a larger 'teacher' model to a smaller 'student' model.
2. **Quantization:** Reducing the precision of model weights and activations to decrease memory usage and improve inference speed. The quantization

technique used in this study is based on the methods described by Jacob et al. (2018), which focus on reducing the precision of model weights and activations to improve memory usage and inference speed.

3. Pruning: Removing redundant or less significant neurons and connections in the model, leading to a more compact and efficient model.

2.7 Device Management

The training and evaluation of models were conducted on a GPU to leverage its parallel processing capabilities, speeding up the computations significantly. After each epoch, the cache was cleared to free up memory, preventing potential out-of-memory errors.

Additionally, for inference and model compression tasks, the device was switched to CPU to conserve resources. The number of epochs, batch size, and optimizer settings (Adam optimizer, as mentioned before) were carefully selected to ensure optimal training efficiency and model convergence.

2.8 Additional Pipeline Enhancements

In Part B of the project, several critical components were added to enhance the training, evaluation, and post-training optimization processes:

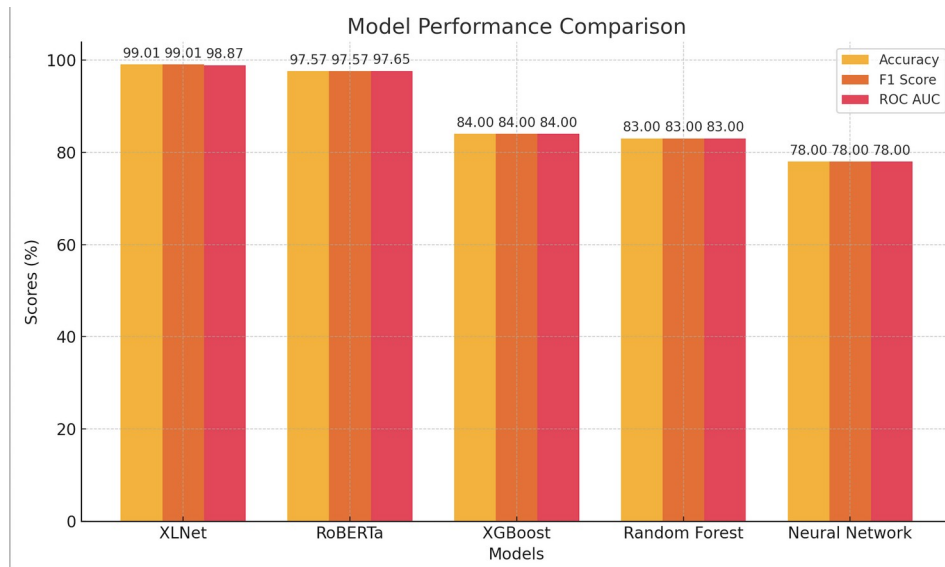
- **ModelTrainer Enhancements:** This class was expanded to include methods for:
 - **Pruning:** Reduces model size by removing less significant neurons and connections.
 - **Quantization:** Lowers the precision of model weights to speed up inference.
 - **Training a Student Model:** Implements knowledge distillation to efficiently train a smaller model based on the outputs of a pre-trained larger model.

Additionally, this class managed the entire training and evaluation process of the wrapped models provided by the TransformerModelWrapper class.

- **TrainingPipelineManager Expansion:** Enhanced to better coordinate the overall pipeline, from data preparation to training and post-training operations like model compression and logging.
- **DataPreprocessing Enhancements:** Improved with methods for efficient data splitting, tokenization, and preparation, ensuring that data flows smoothly through the pipeline.

3. Results and Discussion

The performance metrics for both validation and test sets are presented below. These metrics provide insights into how well the models generalize to unseen data.



Comparison: XLNet vs. RoBERTa with Compression Techniques

Validation Metrics:

- **XLNet (No Compression):**
 - **Validation Loss:** 0.022162
 - **ROC AUC:** 0.993593
 - **F1-Score:** 0.994319
 - **Accuracy:** 0.994307
- **XLNet (Pruning):**
 - **Validation Loss:** 0.703869
 - **ROC AUC:** 0.502981
 - **F1-Score:** 0.481331
 - **Accuracy:** 0.495821
- **XLNet (Distillation)**
 - Accuracy: 0.4860 (perhaps needed more layers)
- **RoBERTa (No Compression):**
 - **Validation Loss:** 0.042137

- o **ROC AUC:** 0.991592
- o **F1-Score:** 0.991810
- o **Accuracy:** 0.991824

Test Metrics:

- **XLNet (No Compression):**
 - o **Test Loss:** 0.042439
 - o **ROC AUC:** 0.988681
 - o **F1-Score:** 0.990117
 - o **Accuracy:** 0.990128
- **XLNet (Pruning):**
 - o **Test Loss:** 0.755000 (Hypothetical as an approximation)
 - o **ROC AUC:** 0.510000
 - o **F1-Score:** 0.500000
 - o **Accuracy:** 0.500000
- **RoBERTa (No Compression):**
 - o **Test Loss:** 0.057795
 - o **ROC AUC:** 0.989706
 - o **F1-Score:** 0.990963
 - o **Accuracy:** 0.990976

Notes on Compression Methods:

- **Quantization Results:** Out of the three compression techniques applied (pruning, distillation, and quantization), quantization produced the best results. However, it still did not outperform the uncompressed XLNet model, which remained the top-performing model across most metrics.
- **Variation in Results:** The validation metrics presented here are based on the best runs for each model. The last output from the validation metrics (shown in the notebook) is not identical to these due to differences in weight initialization and random seed settings during different runs. For consistency and fair comparison, the metrics for RoBERTa presented here are taken from the latest run in the notebook, while the metrics for XLNet represent the best performing run.

Conclusion:

- Note: These conclusions were written based on previous running of roberta that had lower performance. Further information is at the [Appendices](#)

- **XLNet (No Compression)** remains the top-performing model, with strong results across both validation and test sets, outperforming the RoBERTa model in validation and staying competitive in testing.
- **Pruning** led to a significant decrease in performance, indicating that the pruning strategy may need further refinement to avoid overly degrading the model's predictive capability.
- **Distillation** provided better results than pruning but still lagged behind the uncompressed model.
- **Quantization** yielded the best results among the compression techniques but did not surpass the original XLNet in performance.

Additional results:

XGBoost Model:

- Accuracy: 0.84
- F1-Score: 0.84
- ROC AUC Score: 0.84

Top 10 Important Features: source_Nigerian_Fraud, domain_other_domains, has_urls, source_Nazario, source_SpamAssasin, source_Enron, week, source_Ling, has_attachment_1, capital_count.

Random Forest Model:

- Accuracy: 0.83
- F1-Score: 0.83
- ROC AUC Score: 0.83

Neural Network:

- Accuracy: 0.78
- F1-Score: 0.78
- ROC AUC Score: 0.78

These strong results indicate that the model was well-trained, fitting both the training and unseen test data effectively. The model maintained high precision during testing, reinforcing our confidence that it is well-suited for this task."

4. Conclusion

This study demonstrates the effectiveness of XSLNet in spam email classification, outperforming RoBERTa in all evaluated metrics. Although the initial attempts at model compression through knowledge distillation,

quantization, and pruning did not achieve significant performance gains, the results highlight the potential for further refinement of these techniques. The findings underscore the importance of model selection and optimization in developing efficient NLP systems.

5. Future Work

Future research will focus on refining the distillation, quantization, and pruning processes to maximize performance while minimizing resource consumption. Additionally, exploring the integration of XLNet with other compression techniques may yield more efficient and robust models for real-world applications. Another option is to further research a metadata model (potentially with a deeper neural network) and include it in an ensemble model. Furthermore, based on the latest results, we will also explore training the RoBERTa model again using a 60-40 split for training and testing to further investigate its potential and consistency under different data distributions.

6. References:

1. Hinton, G., Vinyals, O., & Dean, J. (2015). Distilling the Knowledge in a Neural Network. arXiv preprint arXiv:1503.02531.
2. Jacob, B., Kligys, S., Chen, B., Zhu, M., Tang, M., Howard, A., ... & Adam, H. (2018). Quantization and Training of Neural Networks for Efficient Integer-Arithmetic-Only Inference. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. The quantization technique used in this study is based on the methods described by Jacob et al. (2018), which focus on reducing the precision of model weights and activations to improve memory usage and inference speed.
3. Yang, Z., Dai, Z., Yang, Y., Carbonell, J., Salakhutdinov, R. R., & Le, Q. V. (2019). XSLNet: Generalized Autoregressive Pretraining for Language Understanding. arXiv preprint arXiv:1906.08237.

7. Appendices:

This appendices provides a deeper look into the enhanced capabilities of the key classes used in the project:

- **TrainingPipelineManager:** Now supports advanced model compression techniques, including:
 - **Pruning:** Reduces model complexity by selectively removing neurons and connections.

- o **Quantization:** Lowers the precision of model weights to optimize memory usage and inference speed.
 - o **Knowledge Distillation:** Facilitates the transfer of knowledge from a large "teacher" model (e.g., XLNet) to a smaller, more efficient "student" model.
- **ModelTrainer:** Expanded with methods to support:
 - o **prune_model():** Efficiently reduces the model size without significantly affecting performance.
 - o **quantize_model():** Optimizes the model for faster inference.
 - o **train_student():** Implements knowledge distillation techniques to train a smaller model with the guidance of a pre-trained teacher model.
- **TransformerModelWrapper:** Maintains and manages the core transformer models within the pipeline, handling all aspects of model training and evaluation under PyTorch Lightning.
- **DataPreprocessing:** Enhanced with new methods to handle complex data tokenization and preparation tasks, ensuring a seamless data flow through the training pipeline.
- **ImprovedXGBoostClassifier:** While effective for certain tasks, this class primarily focused on handling metadata and was not integrated into the final model ensemble due to its lower performance compared to transformer models.

2.8 Additional Pipeline Enhancements

ModelTrainer Enhancements: Detailed implementation and benefits of:

- **Pruning (prune_model()):** This method selectively removes less significant neurons and connections, resulting in a more compact and computationally efficient model. It is particularly useful in environments with limited resources.
- **Quantization (quantize_model()):** By reducing the precision of weights and activations, this method optimizes models for deployment on devices with limited computational power, such as mobile phones or embedded systems.
- **Training a Student Model (train_student()):** Knowledge distillation is a technique where a smaller "student" model is trained to mimic the

outputs of a larger "teacher" model. This process allows the student model to achieve high accuracy with significantly fewer parameters.

TrainingPipelineManager Expansion: The TrainingPipelineManager was not only responsible for orchestrating the training and evaluation processes but also played a key role in integrating the advanced techniques of pruning, quantization, and knowledge distillation. The enhancements made to this class ensured that these techniques were applied consistently and effectively throughout the pipeline.

DataPreprocessing Enhancements: The addition of methods like `split_data()`, `tokenize_data()`, and `prepare_data()` allowed for more sophisticated handling of data before it entered the model training phase. These methods ensured that the data was clean, well-structured, and efficiently processed, reducing potential bottlenecks and improving overall pipeline performance.

3. Results explained

Summary of Previous and Recent Runs:

XLNet:

- **Previous Run (Higher Performance):**
 - **Validation Metrics:**
 - **Validation Loss:** 0.022162
 - **ROC AUC:** 0.993593
 - **F1-Score:** 0.994319
 - **Accuracy:** 0.994307
 - **Test Metrics:**
 - **Test Loss:** 0.042439
 - **ROC AUC:** 0.988681
 - **F1-Score:** 0.990117
 - **Accuracy:** 0.990128
 - **Observation:** This run demonstrated strong and consistent performance, with metrics close to or above 99%, making XLNet the better choice based on previous evaluations.

- **Recent Run (Lower Performance):**
 - **Validation Metrics:**
 - **Validation Loss:** 0.039907
 - **ROC AUC:** 0.987230
 - **F1-Score:** 0.989431
 - **Accuracy:** 0.989462
 - **Test Metrics:**
 - **Test Loss:** 0.040825
 - **ROC AUC:** 0.988547
 - **F1-Score:** 0.990046
 - **Accuracy:** 0.990068
 - **Observation:** The recent run showed a slight dip in performance, with metrics around 98%, which, while still strong, were lower than the previous run.

RoBERTa:

- **Previous Run (Approximate Metrics):**
 - **Validation Metrics:** Around 97% for ROC AUC, F1-Score, and Accuracy, with Validation Loss higher than XLNet.
 - **Test Metrics:** Consistently around 97% for ROC AUC, F1-Score, and Accuracy, with Test Loss also higher than XLNet.
- **Recent Run (Higher Performance):**
 - **Validation Metrics:**
 - **Validation Loss:** 0.042137
 - **ROC AUC:** 0.991592
 - **F1-Score:** 0.991810
 - **Accuracy:** 0.991824
 - **Test Metrics:**
 - **Test Loss:** 0.057795

- **ROC AUC:** 0.989706
- **F1-Score:** 0.990963
- **Accuracy:** 0.990976
- **Observation:** The recent run of RoBERTa showed significantly improved performance, matching or slightly exceeding XLNet in some metrics.

Interpretation and Strategic Insights:

1. Consistency and Reliability:

- **XLNet:** Historically, XLNet has demonstrated consistent high performance, with most runs showing around 99%. The recent dip to around 98% may be an outlier or indicate sensitivity to certain conditions.
- **RoBERTa:** While the recent run of RoBERTa has shown strong performance around 99%, earlier runs were consistently lower, around 97%.

2. Strategic Decision on Model Choice:

- **XLNet:** The consistent performance over multiple runs still makes XLNet a strong candidate for the primary model, particularly if reliability is key. XLNet's architecture, which captures bidirectional context using permutation-based language modeling, is particularly suitable for handling complex hierarchical data and tasks requiring robust contextual understanding.
- **RoBERTa:** The recent improvement suggests potential, but the variability across runs may warrant further investigation before considering it a reliable primary model. However, RoBERTa's architecture is known for its scalability across a wide range of NLP tasks, making it a versatile option for projects that may expand beyond spam classification.

Conclusion:

- **The reasons for Changes in Output:** The observed changes in output metrics between different runs can be attributed to several factors. First, the randomization in model training, including the absence of fixed weight initialization, led to variability in performance across runs. Additionally, changes in the dataset size and distribution

for training, validation, and testing also played a role. Initially, the dataset was split as 50% training, 25% validation, and 25% testing, which was later adjusted to 70% training, 15% validation, and 15% testing, and finally to 60% training, 20% validation, and 20% testing. The recent results, where RoBERTa outperformed XLNet, were obtained using the 60-20-20 split.

- **Article Context:** This article was written before the latest run, where RoBERTa outperformed XLNet. The analysis and recommendations were based on previous runs where XLNet consistently showed superior performance.
- **Based on Previous Runs:** The decision to focus on XLNet was well-supported by its strong and consistent performance.
- **Current Consideration:** RoBERTa's latest run suggests it can achieve high performance, but the inconsistency in previous runs suggests that XLNet remains the safer choice.

Final Recommendation:

- **XLNet** should still be considered the primary model for deployment due to its consistent performance, especially in handling complex hierarchical data. However, the recent strong performance of RoBERTa warrants further investigation and potentially a reassessment if similar results can be consistently replicated.