

# Abstract

This report presents a streamlined narrative of a deep-learning workflow for fine-tuning a transformer-based text classifier on the AG News dataset. Leveraging a pre-trained DistilBERT model, the pipeline encompasses dataset acquisition, text encoding, lightweight prototyping via dataset subsampling, model adaptation for four-way classification, and evaluation under standardized metrics. Key design choices, including transfer-learning strategies, dynamic padding, and hyperparameter selection, are elaborated to convey both conceptual underpinnings and practical considerations.

## Introduction

Text classification remains a core task in natural language processing, with applications ranging from news categorization to sentiment analysis. Recent advances in Transformer architectures—particularly BERT and its distilled variants—have demonstrated state-of-the-art performance across diverse benchmarks. This project harnesses such a model, DistilBERT, pre-trained on a sentiment-analysis corpus, and adapts it to the multi-class AG News task via fine-tuning. The report refrains from code listings, instead emphasizing conceptual workflow and rationale.

## Data Acquisition and Preprocessing

The AG News corpus consists of four topical categories (World, Sports, Business, Sci/Tech), each with tens of thousands of examples. The workflow begins by loading this dataset through a high-level API, which automatically retrieves and partitions the data into training and test splits. To prepare text for model ingestion:

**Tokenization:** A DistilBERT tokenizer converts raw sentences into fixed-length token sequences, applying both truncation (to manage overly long texts) and padding (to equalize batch dimensions).

**Column Management:** Non-essential fields (e.g., original text) are stripped post-tokenization, and label columns are reformatted to align with the model's expected input.

**Tensor Formatting:** The processed data is cast into tensor form, enabling seamless integration with the deep-learning framework's DataLoader utilities.

To facilitate rapid experimentation, a small subset (1,000 samples each for train and eval) is randomly selected, ensuring reproducibility via a fixed seed. This lightweight prototyping stage accelerates hyperparameter tuning without sacrificing representativeness.

## Model Adaptation and Architecture

At the core of the pipeline lies DistilBERT—a compact, faster variant of BERT that retains the bulk of its representational power at reduced computational cost. Two modifications enable its transition from binary sentiment analysis to four-way news classification:

**Classification Head Replacement:** The original two-class output layer is substituted with a new linear head sized to four outputs.

**Checkpoint Loading:** To maximize reuse of pre-trained weights, the Transformer backbone is initialized from the sentiment-analysis checkpoint, while the classification head is randomly initialized, thanks to a loading flag that tolerates size mismatches.

This strategy embodies transfer learning, where foundational language representations are co-opted for a novel downstream task, reducing both training time and data requirements.

## Training Configuration

Fine-tuning is orchestrated via a high-level training API that encapsulates optimization, scheduling, and checkpointing. Key settings include:

**Learning Rate:** A modest initial rate ( $3 \times 10^{-5}$ ) balances convergence speed with stability on a small dataset.

**Batch Sizes:** Training and evaluation batches are sized to fit GPU memory constraints while ensuring statistical robustness (e.g., 16 for train, 32 for eval).

**Epochs:** Three full passes through the data are specified, offering sufficient gradient updates without overfitting.

**Weight Decay:** A small regularization term (0.01) discourages over-complex weight patterns, promoting generalization.

**Evaluation Strategy:** Performance is assessed at the end of each epoch, with model checkpoints saved correspondingly.

A dynamic padding collator further refines efficiency by padding only to the maximum sequence length present within each batch, reducing wasted computation on excess padding tokens.

## Evaluation Metrics and Results

Model performance is evaluated on the held-out subset using standard classification metrics:

**Accuracy:** Proportion of correctly predicted labels across all classes.

**Precision & Recall:** Computed per class and macro-averaged to gauge both exactness and completeness.

**F1-Score:** Harmonic mean of precision and recall, offering a balanced measure under class imbalance.

On the 1,000-example test subset, fine-tuning yields competitive accuracy and balanced precision/recall trade-offs, demonstrating the viability of a distilled Transformer even under constrained data. Detailed metric logs across epochs provide insight into learning dynamics and convergence behavior.

## Discussion and Recommendations

While the distilled Transformer adapts effectively, several avenues exist for further enhancement:

**Full-Dataset Training:** Scaling from the prototyping subset to the full AG News corpus would likely boost performance and resilience to rare category examples.

**Advanced Tokenization Schemes:** Experimentation with longer maximum lengths or alternative tokenizers (e.g., RoBERTa's Byte-Level BPE) could capture additional contextual nuances.

**Learning-Rate Schedules:** Incorporating warmup phases or cosine-decay schedules may improve convergence stability.

**Data Augmentation:** Techniques such as back-translation or synonym replacement can enrich the training set, reducing over-dependence on limited textual patterns.

**Error Analysis:** Confusion-matrix inspection and per-class performance breakdowns would pinpoint systematic misclassifications, guiding targeted refinements.

**Cross-Validation:** K-fold splits across the training set would provide more robust generalization estimates and mitigate split-specific artifacts.

## Conclusion

This pipeline exemplifies a concise yet powerful approach to repurposing a distilled Transformer for multi-class text classification. By combining transfer learning, dynamic batching, and a succinct training regimen, it achieves prompt convergence and strong baseline performance. Future work should expand scope to the full dataset, incorporate richer evaluation analyses, and explore complementary architectures for even greater accuracy and robustness.

## References

- Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding.
- Sanh, V., Wolf, T., & Rush, A. M. (2019). DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter.
- Zhang, X., Zhao, J., & LeCun, Y. (2015). Character-level Convolutional Networks for Text Classification.
- “AG News” Dataset. Hugging Face Datasets.