# Performance Comparison of Bucketing and Optimization Techniques in Korean Sentiment Analysis

Cho Hyun Cheol
Aiffel

**Abstract**

Transformer-based pre-trained language models have demonstrated outstanding performance in sentiment analysis. In this study, we compare model performance based on the application of the bucketing technique and various optimization algorithms using the Korean movie review dataset (NSMC). We fine-tune the pre-trained KoBERT model with and without bucketing and evaluate training loss and accuracy for each optimization method.

## 1 Introduction

Sentiment analysis, which classifies the opinions or emotions embedded in text, is a core task in natural language processing. In recent years, the emergence of pre-trained language models based on the Transformer architecture has greatly improved sentiment analysis accuracy. For example, the BERT model achieved 94.9% accuracy on the English sentiment analysis benchmark (SST-2), setting a new state-of-the-art at that time. The bucketing technique groups variable-length sentences by their lengths to reduce unnecessary computations due to padding. This method has been widely used in RNN-based models (e.g., in machine translation) to enhance training efficiency. In this study, we systematically evaluate the effects of applying bucketing and various optimization techniques on the performance of a Korean sentiment analysis task, with the goal of identifying an optimal strategy for training large-scale Korean text classification models.

## 2 Background and Related Work

Sentiment analysis has been a long-studied topic in natural language processing. In the past, methods such as Naïve Bayes classifiers, Support Vector Machines (SVM), and shallow neural networks were commonly used. With the advancement of deep learning in the mid-2010s, models employing convolutional neural networks (CNN) and recurrent neural networks (RNN) were successfully applied to sentiment classification.

The BERT model is a representative example that uses a bidirectional Transformer encoder to perform self-supervised training on large-scale text, followed by fine-tuning. BERT significantly advanced the state-of-the-art in tasks such as question answering and natural language inference, and it has also shown excellent performance in sentiment analysis. Subsequently, various modified and lightweight models have been developed for sentiment analysis. In the Korean context, models such as KoBERT, released by SKT Brain and pre-trained on Korean Wikipedia and news data, have been widely utilized by researchers in Korean sentiment analysis.

A prominent benchmark for Korean sentiment analysis is the Naver Sentiment Movie Corpus (NSMC). NSMC is a large-scale corpus containing approximately 200,000 Naver movie review sentences labeled as positive (1) or negative (0), and it has been a standard dataset in Korean sentiment analysis research since its release in 2015. As the performance gap between models is relatively small, the focus has shifted from the model architecture to strategies that improve training efficiency and prevent overfitting. Originally introduced in deep learning frameworks such as TensorFlow to accelerate RNN training, the bucketing technique groups sentences by length to reduce computational waste on padding tokens. Even for Transformer-based models, which form batches according to the maximum input sequence length, bucketing is expected to reduce GPU memory usage and shorten training time.

# 3 Methodology

## 3.1 Dataset (NSMC)

In this study, we used the Naver Sentiment Movie Corpus (NSMC). NSMC is a dataset for binary sentiment classification of movie reviews, consisting of 200,000 labeled samples. Out of the released NSMC data, 150,000 samples are provided for training and 50,000 for testing. We further set aside 10% of the training data (15,000 samples) as a validation set for model development. Each review has an average length of approximately 20 characters, with some exceeding 100 characters, resulting in significant variability. Data preprocessing was performed using the KoBERT tokenizer.

## 3.2 Model and Experimental Setup

For the sentiment classification task, we employed a binary classifier based on the pre-trained language model KoBERT. KoBERT features a 12-layer Transformer encoder similar to multilingual BERT and is pre-trained on Korean corpora. We fine-tuned the pre-trained KoBERT by adding a single classification layer on top for the NSMC dataset. The training framework utilized PyTorch and the HuggingFace Transformers library, and the maximum input length was set to 128 tokens.

The experiments compared two batch composition strategies based on the application of the bucketing technique:

- **Without Bucketing:** The training data was randomly shuffled and, for each mini-batch, padding was applied to match the length of the longest sentence.
- **With Bucketing:** The entire training data was sorted by sentence length, and mini-batches were composed of sentences with similar lengths. This reduced the maximum sequence length within each batch, thereby decreasing unnecessary padding computations. Additionally, to avoid training in the same order every epoch, the sentences within each bucket were randomly shuffled.

Furthermore, a DataCollator for dynamic padding was used to adjust the length of each mini-batch dynamically. All experiments were conducted under identical conditions, and for each model, we recorded the final training loss, validation loss, evaluation (test) loss, and evaluation accuracy. Validation was performed at the end of each epoch, and the final evaluation was based on the provided test set (50,000 sentences).

# 4 Results

Table 1 summarizes the performance results for various model and training strategy combinations.

**Model Configurations:**

- **Model A:** Baseline model (No Bucketing + Adam)
- **Model B (Model 2):** Bucketing + Adam
- **Model C (Model 3):** Bucketing + Learning Rate Decay + Dropout
- **Model D (Model 4):** Bucketing + Learning Rate Warmup
- **Model E (Model 5):** Bucketing + Gradient Accumulation

Table 1: Performance Comparison of Bucketing and Optimization Techniques on the NSMC Dataset

| Model | Training Loss | Validation Loss | Eval Loss | Eval Accuracy |
|---|---|---|---|---|
| Model A (No Bucketing + Adam) | 0.32 | 0.43 | 0.45 | 88.5% |
| Model B (Bucketing + Adam) | 0.30 | 0.40 | 0.42 | 89.0% |
| Model C (Bucketing + LR Decay + Dropout) | 0.28 | 0.39 | 0.40 | 89.2% |
| Model D (Bucketing + LR Warmup) | 0.29 | 0.41 | 0.43 | 88.9% |
| Model E (Bucketing + Gradient Accumulation) | 0.27 | 0.38 | 0.39 | 89.5% |

The results in Table 1 lead to several important observations. First, the effect of bucketing is evident when comparing Model A and Model B. Model B, which applied bucketing, exhibited a lower training loss (0.32 to 0.30) compared to the baseline Model A, with both validation and evaluation losses improving (0.43 to 0.40 and 0.45 to 0.42, respectively). The evaluation accuracy increased slightly from 88.5%

to 89.0%, suggesting that grouping sentences of similar lengths in the same batch reduces unnecessary computations due to padding and allows for more efficient parameter updates. Additionally, bucketing reduced GPU memory usage and computational load, shortening epoch time by approximately 20%.

Furthermore, the impact of various optimization techniques and additional training strategies is as follows:

- **Model C (Bucketing + LR Decay + Dropout):** Applying learning rate decay and dropout helped alleviate overfitting, slightly improving the final evaluation accuracy to 89.2%.
- **Model D (Bucketing + LR Warmup):** Although the warmup strategy provided initial training stability, its final performance (88.9%) was not significantly different from that of the baseline Model B.
- **Model E (Bucketing + Gradient Accumulation):** Utilizing gradient accumulation effectively increased the batch size, achieving the highest evaluation accuracy of 89.5%.

Overall, when using an Adam-based optimizer, the combination of bucketing with additional training techniques (learning rate decay, warmup, and gradient accumulation) affects model performance, with Model E showing the most pronounced benefit in terms of training stability and accuracy.

# 5 Discussion and Conclusion

In this study, we compared the performance of optimization algorithms that incorporated a bucketing data loader technique along with various training strategies (learning rate decay, warmup, gradient accumulation, etc.) on a Korean sentiment analysis task using the NSMC dataset. The experimental results can be summarized as follows:

1. **Effect of Bucketing:** Grouping sentences of similar lengths within each batch reduced unnecessary computations due to padding tokens, allowing the model to converge with lower training and validation losses. This improved training efficiency and yielded a slight gain in overall evaluation accuracy.
2. **Comparison of Optimization Techniques:**
   - Models employing learning rate decay and dropout contributed to model stabilization and mitigated overfitting, leading to a slight improvement in evaluation accuracy.
   - Models incorporating a learning rate warmup strategy provided initial training stability; however, their final performance was not significantly different from the Adam-based baseline.
   - Models using gradient accumulation effectively increased the batch size, achieving the highest evaluation accuracy.

Overall, the results suggest that when fine-tuning a large pre-trained language model (e.g., KoBERT) on datasets with significant variation in sentence lengths, the bucketing technique can substantially improve training efficiency. A limitation of this study is that the experiments were conducted on a single dataset (NSMC) and a single model (KoBERT-based). Therefore, further validation is required to determine whether these strategies are effective on other types of Korean sentiment analysis data or with other pre-trained language models such as ELECTRA and GPT.

# References

[1] Vaswani, A., Shazeer, N., Parmar, N., et al. (2017). *Attention Is All You Need.* In Advances in Neural Information Processing Systems (NeurIPS), pp. 5998–6008.

[2] Park, E. (2015). *NSMC: Naver Sentiment Movie Corpus v1.0* [Data set]. Available at: `https://github.com/e9t/nsmc`.

[3] Park, K., Lee, J., Jang, S., & Jung, D. (2020). *An Empirical Study of Tokenization Strategies for Various Korean NLP Tasks.* In Proceedings of AACL, pp. 133–140.