

# Optimizing Sentiment Classification using DSPy's MIPRO Optimizer

## Introduction

This report documents the application of prompt engineering using the DSPy package for a sentiment analysis task. Initially, I planned to use the OLLAMA model, but due to hardware limitations, I switched to OpenAI's GPT-3.5-turbo. The task involves classifying IMDB movie reviews as either positive or negative.

## Dataset Preparation

I used the IMDB dataset, a widely-used benchmark for sentiment analysis. However, due to OpenAI's API rate limits, I worked with a subset of 200 reviews instead of the full dataset of 50,000. This ensured stable API access throughout the experiment.

## Problem Definition

The task is binary sentiment classification, where:

- Input: A movie review text.
- Output: A sentiment label (positive/negative).

This aligns with the DSPy LM program signature, which maps text input to sentiment labels.

## Evaluation Metric

I used accuracy as the primary evaluation metric. It measures how often the model correctly classifies the sentiment of a given review. This provides a straightforward assessment of the model's performance.

## Baseline Model Performance

The baseline model, using DSPy's default settings, achieved 98% accuracy on the test set. This established a strong benchmark before optimization.

## Optimized Model Performance

I applied DSPy's MIPROv2 optimizer to refine the prompt structure and enhance classification robustness. The optimized model achieved 93.8% accuracy, reflecting a slight decrease from the baseline.

### Why Did Accuracy Drop?

- The MIPROv2 optimizer improved prompt clarity and instruction alignment, leading to more detailed sentiment analysis.
- The model generated more nuanced outputs, which may have introduced minor inconsistencies in classification.
- Rather than a straightforward classification, the optimized model provided richer reasoning, potentially impacting classification strictness.

More details, including specific error cases and evaluation insights, can be found in the Jupyter notebook's evaluation results.

## Optimization Process

The MIPROv2 optimizer improved the model in three key steps:

### 1. Bootstrapping Few-Shot Examples

- The optimizer generated few-shot learning examples to enhance model generalization.

Example Log (from the notebook):

plaintext

CopyEdit

2025/03/11 11:25:03 INFO dspy.teleprompt.mipro\_optimizer\_v2:

==> STEP 1: BOOTSTRAP FEWSHOT EXAMPLES <==

### 2. Instruction Proposal & Refinement

- The optimizer suggested multiple prompt variations, evaluating which structure performed best.

Example Proposed Instruction:

"Analyze the sentiment of movie reviews and provide reasoning for the sentiment classification."

### 3. Iterative Prompt Optimization

- DSPy iteratively evaluated different combinations of prompts and few-shot examples to refine prompt effectiveness.
- *Note:* While some prompt optimization techniques use Bayesian optimization, DSPy primarily tests variations iteratively to identify the best-performing configuration.

#### Before-and-After Prompt Examples

Below is an example of how prompt instructions evolved after optimization:

##### Baseline Prompt:

- *"Classify the sentiment of this movie review as 'positive' or 'negative'."*

##### Optimized Prompt (Generated by DSPy's MIPROv2 optimizer):

- *"Analyze the sentiment of the following movie review. Identify words that indicate a positive or negative sentiment and provide reasoning for the classification. Respond with 'positive' or 'negative' only."*

The refined prompt encouraged deeper analysis, making model outputs more explanatory. More detailed prompt evolution can be found in the Jupyter notebook section on prompt refinement.

#### Insights and Comments

- The optimized prompts made the model more interpretable but slightly less precise in strict classification.
- While accuracy dropped, the interpretability and reasoning depth of responses significantly improved.
- The API rate limit issues constrained the dataset to 200 samples, which may have introduced data variability affecting final accuracy.
- Switching from OLLAMA to GPT-3.5-turbo allowed stable experimentation despite hardware constraints.

#### API Call Log Example

The API calls were executed using DSPy's MIPROv2 optimizer, which iteratively refined prompts to enhance classification performance. An example API call log can be found in the Jupyter notebook, showing execution details such as timestamp and API response.

## Example API Call Log (from the notebook):

2025/03/11 11:26:17 INFO dspy.evaluate.evaluate:  
Average Metric: 75 / 80 (93.8%)

More detailed logs are documented in the submitted Jupyter notebook.

## Conclusion

This assignment demonstrated DSPy's effectiveness in prompt optimization for sentiment classification. The MIPROv2 optimizer improved prompt clarity and model interpretability, though at a slight cost in accuracy.

## Future Improvements

- Experimenting with different DSPy optimizers.
- Adjusting few-shot learning examples to balance accuracy vs. reasoning depth.
- Using larger datasets if API constraints allow.

## References

1. DSPy Tutorials – Official guide on entity extraction. Available at: [https://dspy.ai/tutorials/entity\\_extraction/](https://dspy.ai/tutorials/entity_extraction/)
2. OpenAI API Documentation – Guide on GPT-3.5-turbo usage. Available at: <https://platform.openai.com/docs/>
3. YouTube Resources – Tutorials on DSPy optimization strategies.
4. ChatGPT & Perplexity AI – Used for debugging and refining DSPy implementation.

## Source Code Availability

The complete implementation, including dataset preparation, model training, DSPy optimization, and evaluation scripts, is available in the submitted Jupyter notebook and also on GitHub: [GitHub Repository](#). The notebook contains API call logs, prompt evolution, model evaluations, and results.