

Overview

1. Seeking to alleviate the constraints of conventional Weight Averaging (WA) in the fine-tuning of Pretrained Language Models (PLMs).
2. Intending to propose an efficient model fusion methodology that incorporates Hyperparameter Optimization (HP).

Contributions

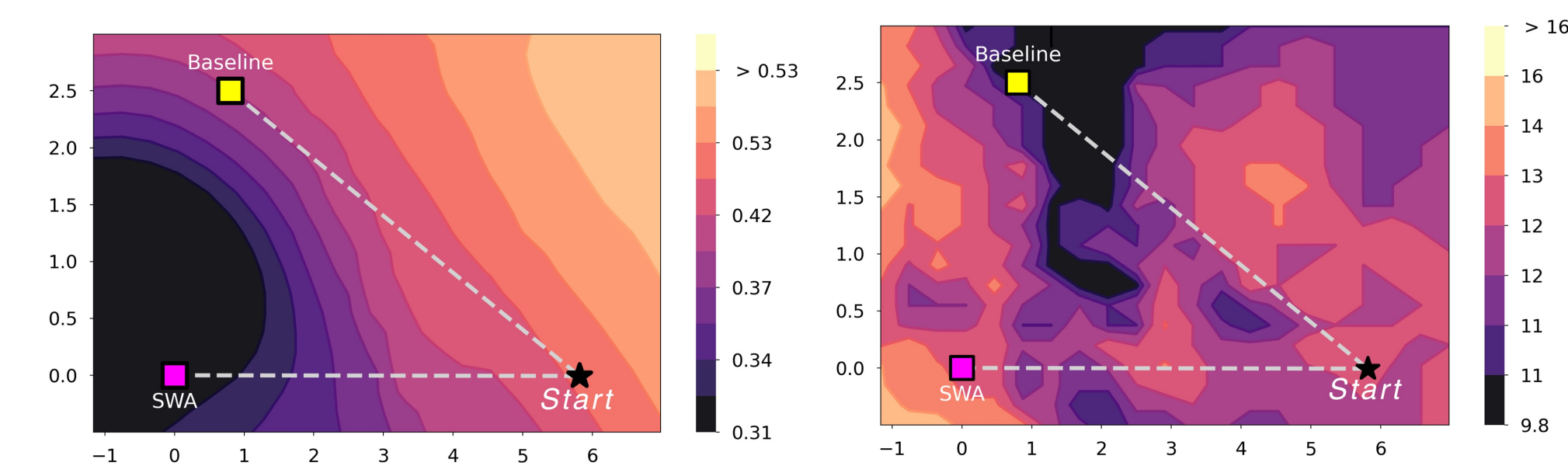
1. **Multi Objective Bayesian Optimization SWA (MOBO-SWA):**
 - Employing MOBO optimizes model averaging coefficients by considering both loss and metric.
2. **Bilevel Model Fusion (Bilevel-BO-SWA):**
 - The process is structured as a bilevel procedure, with an outer BO optimizing hyperparameters for PLMs fine-tuning, and an MOBO-SWA for model fusion.

Backgrounds

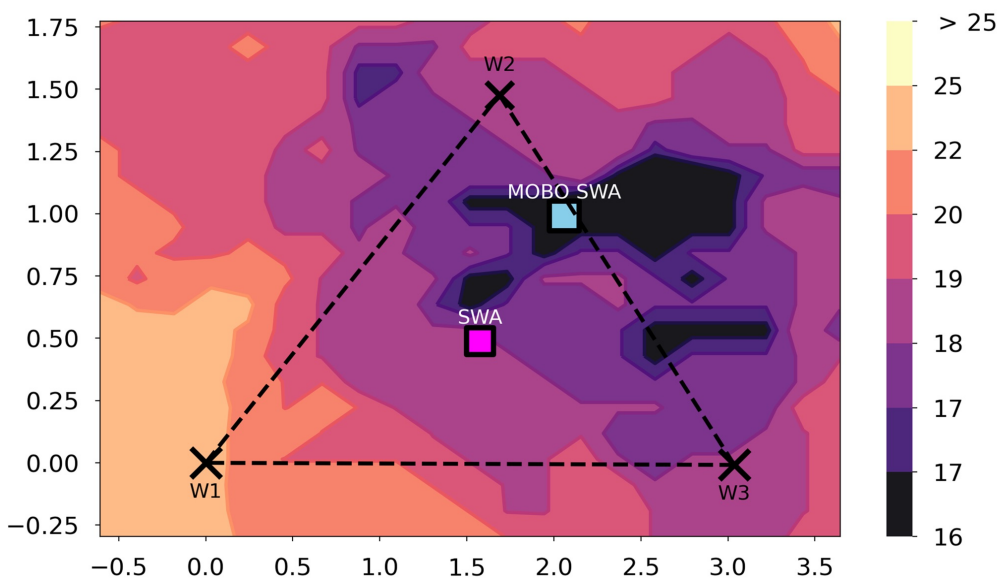
- **Model Fusion for Pre-trained Language Models**
 - The cost of fine-tuning PLMs is **high**, making deep ensemble methods **inefficient**.
 - WA methods like **Stochastic Weight Averaging (SWA)** [1] and **Model Soups** [2] are more feasible for model fusion, although not always adequate for PLMs [3].
 - We discovered that the loss surface of PLMs has a best combination that cannot be found by simple averaging.
- **Bayesian Optimization (BO)**
 - BO is a strategy for optimizing black-box functions that are **costly to evaluate** using the surrogate model.
 - Applications range across various domains, and while Gaussian Processes are commonly used as the surrogate model, other models like Bayesian neural networks and tree-based models are also options [4].

Loss Surface of the Language Model

Valid loss surface vs Valid metric surface



SWA vs MOBO-SWA

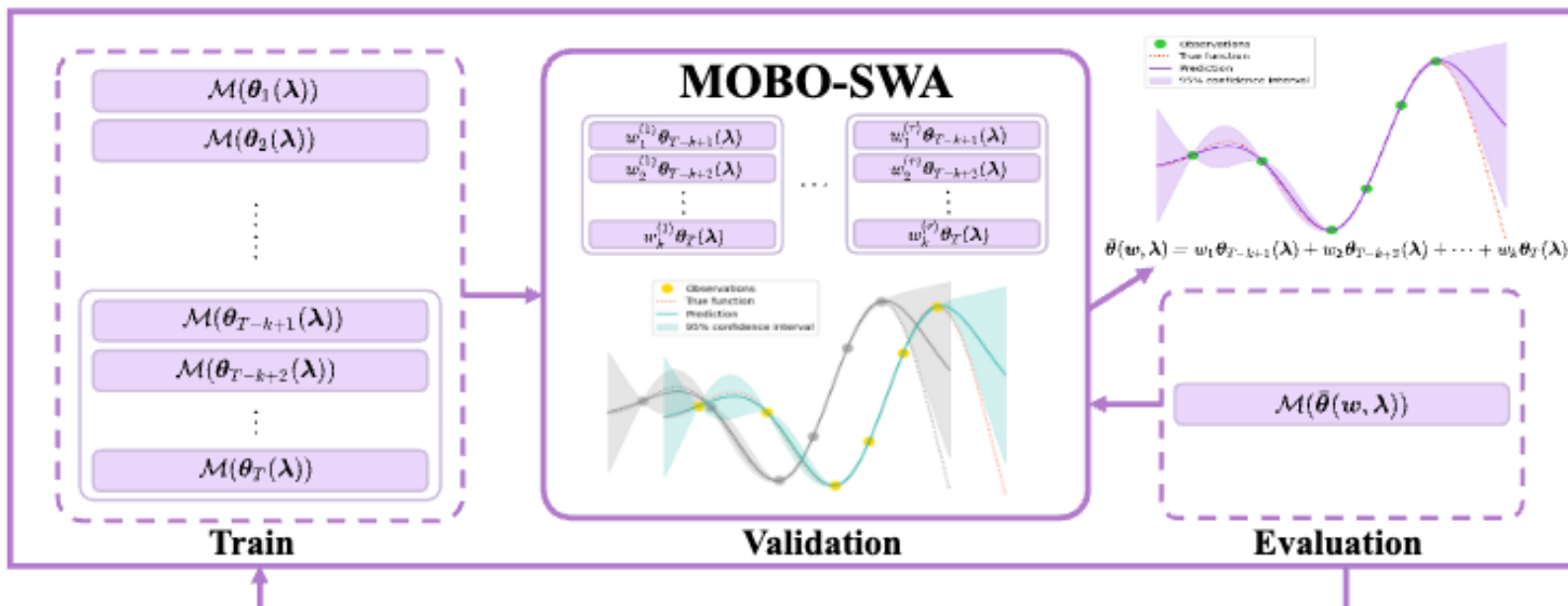


- ✓ While SWA tends to lead to a **flat loss minima**, this is not necessarily related to better performance according to **evaluation metrics**.
- ✓ There exists the **best combination** for combining model weights that is not achievable through SWA.

* Following Garipov et al.'s [5] method, we collect SWA models w_1, w_2, w_3 and establish an *orthonormal basis* u_1, u_2 to represent their parameter space on the x-axis and y-axis.

Model Fusion Through Bayesian Optimization

Bilevel-BO-SWA



- Our target language model as $\mathcal{M}(\theta(\lambda))$ where λ is a hyperparameter vector and $\theta(\lambda)$ is model parameter trained with λ

$$\bar{\theta}(w, \lambda) = w_1 \theta_{T-k+1}(\lambda) + w_2 \theta_{T-k+2}(\lambda) + \dots + w_k \theta_T(\lambda), \quad (1)$$
 Where $w_1, w_2, \dots, w_k \in [0, 1]$ subject to $\sum_{i=1}^k w_i = 1$
- Our objective is to aggregate the final k epoch models into a proficient single model $\mathcal{M}(\bar{\theta}(\lambda))$, using a weighted combination shown in equation (1).

Multi-Objective Bayesian Optimization for Model Fusion

$$\mathcal{P} = \left\{ w^\dagger \mid w^\dagger = \arg \min_w (f_{\text{loss}}(\mathcal{M}(\bar{\theta}(w, \lambda))), f_{\text{metric}}(\mathcal{M}(\bar{\theta}(w, \lambda))) \right\}. \quad (2)$$

We optimize model fusion coefficients with MOBO-SWA (2) considering both loss and metric, employing qNEHVI [6] for the hypervolume improvement objective optimization.

Bilevel Bayesian Optimization for Model Fusion

We further refine hyperparameters through Bilevel-BO-SWA with the objective $f_{\text{metric}}(\mathcal{M}(\bar{\theta}(w^\dagger, \lambda)))$

Experiment Results

Main result

Table 1: **Results on the GLUE dataset using the RoBERTa-base.** Numerical results in boldface and with an underscore indicate the best and the second-best results in the respective datasets, respectively.

| Method | RTE | MRPC | CoLA | STS-B | SST-2 | Avg. |
|-----------------------|-------------------------|-------------------------|-------------------------|-------------------------|-------------------------|--------------|
| Fine-tune | 74.94 \pm 2.28 | 91.65 \pm 0.65 | 56.34 \pm 2.90 | 89.86 \pm 0.16 | 94.49 \pm 0.04 | 81.46 |
| SWA [16] | 77.19 \pm 1.04 | 91.31 \pm 1.82 | 55.05 \pm 3.09 | 89.89 \pm 0.20 | 94.49 \pm 0.08 | 81.59 |
| Greedy SWA [37] | 76.52 \pm 1.31 | 91.84 \pm 0.14 | 56.47 \pm 3.20 | 89.87 \pm 0.18 | 94.36 \pm 0.05 | 81.81 |
| Learned SWA [37] | 77.82 \pm 3.58 | 90.62 \pm 1.87 | <u>59.02</u> \pm 2.60 | 89.65 \pm 0.09 | 94.19 \pm 0.00 | 82.26 |
| MOBO-SWA (ours) | <u>77.86</u> \pm 0.28 | <u>92.05</u> \pm 1.05 | 58.20 \pm 1.72 | 89.58 \pm 0.12 | <u>94.55</u> \pm 0.12 | <u>82.45</u> |
| Bilevel-BO-SWA (ours) | 78.43 \pm 0.26 | 92.38 \pm 0.68 | 59.21 \pm 3.53 | 89.86 \pm 0.01 | 94.97 \pm 0.08 | 82.97 |
| Best subset (oracle) | 80.66 \pm 0.52 | 92.90 \pm 0.22 | 60.01 \pm 1.88 | 89.93 \pm 0.20 | 95.08 \pm 0.06 | 83.72 |

Loss vs Metric discrepancy

Table 2: **Results on GLUE benchmark for RoBERTa-base.** Evaluation results of SWA and naive fine-tuned model on the RTE, MRPC, SST-2. We used custom validation sets for the evaluation. Here *NLL* is the loss function and *Error rate* is the 1 - *Accuracy* for the RTE and SST-2, and the *F1 score* for the MRPC. The lower value is the better for all the evaluation functions. Please refer to **Appendix B** to see how we split the custom validation sets.

| | | Task | | |
|-----------------------------|-----------|-------------------------|------------------------|------------------------|
| | | RTE | MRPC | SST-2 |
| NLL (\downarrow) | Fine-tune | 0.97 \pm 0.01 | 0.54 \pm 0.02 | 0.28 \pm 0.00 |
| | SWA | 0.87 \pm 0.03 | 0.53 \pm 0.00 | 0.22 \pm 0.00 |
| Error rate (\downarrow) | Fine-tune | 21.21 \pm 0.69 | 7.82 \pm 0.01 | 4.94 \pm 0.26 |
| | SWA | 21.71 \pm 1.47 | 7.90 \pm 0.01 | 5.16 \pm 0.24 |

Ablation on MOBO and BO

Table 3: **Using RoBERTa-base, Performance Analysis of Basic BO on the GLUE Dataset.** When employing BO that focuses solely on a single objective, specifically the metric, it was observed that MOBO-SWA exhibited commendable effectiveness in comparison to BO-SWA, which takes into account both the loss and metric.

| Method | RTE | MRPC | CoLA | STS-B | SST-2 | Avg. |
|----------|-------------------------|-------------------------|------------------|-------------------------|-------------------------|--------------|
| BO-SWA | 77.20 \pm 1.97 | 91.92 \pm 0.66 | 57.56 \pm 0.30 | 89.63 \pm 0.05 | 94.47 \pm 0.15 | 82.16 |
| MOBO-SWA | <u>77.86</u> \pm 0.28 | <u>92.05</u> \pm 1.05 | 58.20 \pm 1.72 | 89.58 \pm 0.12 | <u>94.55</u> \pm 0.12 | <u>82.45</u> |

Ablation on the Effectiveness of the Outer BO

Table 4: **Comparative Performance Analysis Applying Outer BO on Various Baselines.** The table shows that Bilevel-BO-SWA outperforms other strategies on RTE and MRPC datasets, according to key performance metrics.

| | Baseline | SWA | Greedy SWA | Learned SWA | Bilevel-BO-SWA |
|------|----------|-------|------------|-------------|----------------|
| RTE | 77.20 | 76.68 | 76.68 | 78.22 | 79.60 |
| MRPC | 91.41 | 90.57 | 90.57 | 90.03 | 93.39 |

References

- [1] Izmailov, Pavel, et al. "Averaging Weights Leads to Wider Optima and Better Generalization."
- [2] Wortsman, Mitchell, et al. "Model soups: averaging weights of multiple fine-tuned models improves accuracy without increasing inference time." International Conference on Machine Learning. PMLR, 2022.
- [3] Kaddour, Jean, et al. "When do flat minima optimizers work?." Advances in Neural Information Processing Systems 35 (2022): 16577-16595.
- [4] Frazier, Peter I. "A tutorial on Bayesian optimization." arXiv preprint arXiv:1807.02811 (2018).
- [5] Garipov, Timur, et al. "Loss surfaces, mode connectivity, and fast ensembling of dnns." Advances in neural information processing systems 31 (2018).
- [6] Daulton, Samuel, Maximilian Balandat, and Eytan Bakshy. "Parallel bayesian optimization of multiple noisy objectives with expected hypervolume improvement." Advances in Neural Information Processing Systems 34 (2021): 2187-2200.