

---

**Algorithm:** Lightweight data collection for MAP regression (inter-leaved pruning).

---

**Input:** Pre-trained model  $Y_0$ , number of rounds  $rds$ , pruning budget  $k$ , total layers  $L$

**Output:** Training dataset  $train\_data$

```
1 Initialize  $train\_data \leftarrow \{(\tilde{m}_0^a, \tilde{m}_0^g, a_0)\}$ ;
2 Function PruneIterative $((\tilde{m}_0^a, \tilde{m}_0^g), (\tilde{m}_{\max}^a, \tilde{m}_{\max}^g))$ :
3   for  $n \leftarrow 1$  to  $rds$  do
4     // Step 1: prune attention
5      $\tilde{m}_n^a \leftarrow \tilde{m}_0^a + n \cdot \frac{\tilde{m}_{\max}^a - \tilde{m}_0^a}{rds};$ 
6      $\tilde{m}_{n-1}^g \leftarrow \tilde{m}_0^g + (n-1) \cdot \frac{\tilde{m}_{\max}^g - \tilde{m}_0^g}{rds};$ 
7     Prune  $Y_{n-1}$  along attention to ratio  $\tilde{m}_n^a$  to obtain  $Y'_n$ ;
8     // Fine-tuning and evaluation
9     Fine-tune  $Y'_n$  and evaluate  $\rightarrow (\tilde{m}_n^a, \tilde{m}_{n-1}^g, a_n)$ ;
10    Append  $(\tilde{m}_n^a, \tilde{m}_{n-1}^g, a_n)$  to  $train\_data$ ;
11    // Step 2: prune activation
12     $\tilde{m}_n^g \leftarrow \tilde{m}_0^g + n \cdot \frac{\tilde{m}_{\max}^g - \tilde{m}_0^g}{rds};$ 
13    Prune  $Y'_n$  along activation to ratio  $\tilde{m}_n^g$  to obtain  $Y_n$ ;
14    // Fine-tuning and evaluation
15    Fine-tune  $Y_n$  and evaluate  $\rightarrow (\tilde{m}_n^a, \tilde{m}_n^g, a_n)$ ;
16    Append  $(\tilde{m}_n^a, \tilde{m}_n^g, a_n)$  to  $train\_data$ ;
17  end
18 Set  $\tilde{m}_{a,\max} \leftarrow k/L$  and  $\tilde{m}_{g,\max} \leftarrow k/L$ ;
19 PruneIterative $((\tilde{m}_0^a, \tilde{m}_0^g), (\tilde{m}_{a,\max}, \tilde{m}_{g,\max}))$ ;
20 return  $train\_data$ ;
```

---