

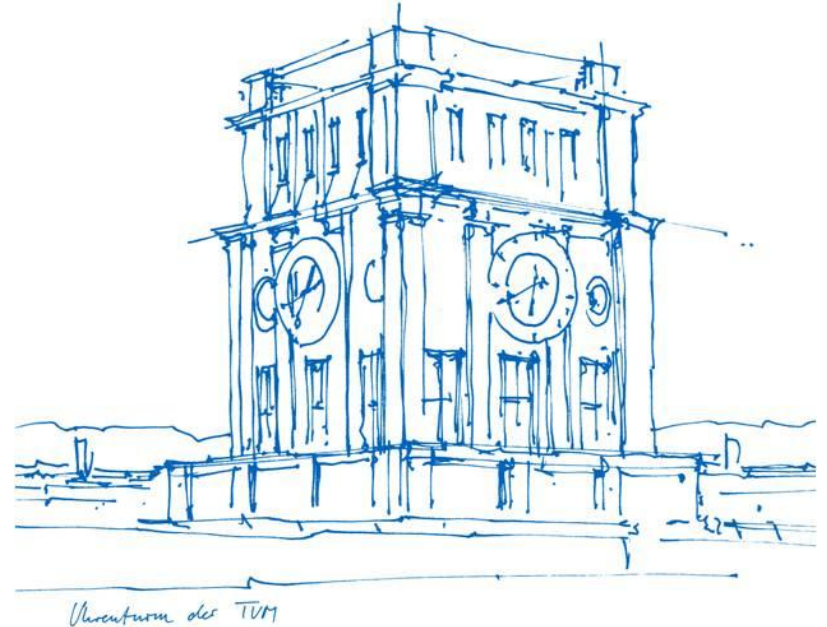
SELECTIVE PARAMETER UPDATING - WEEK 1

Coşku Barış Coşlu

Mato Gudelj

Baykam Say

03 May 2023



Contents

1. Background Research
2. Downstream Task & Dataset
3. Training Loop
4. Generic Update Policy Hook
5. Open Issues & Solutions
6. Plan for the Next Two Weeks

Background Research - Recommended Readings

- Selective Backprop [1]
 - Perform forward pass (every nth epoch)
 - Skip backward pass if loss is too low
- Weight Update Skipping [2]
 - Skip updating weights and only update bias
 - Go back to normal training if no improvements occur
 - Only update biases for the last few layers

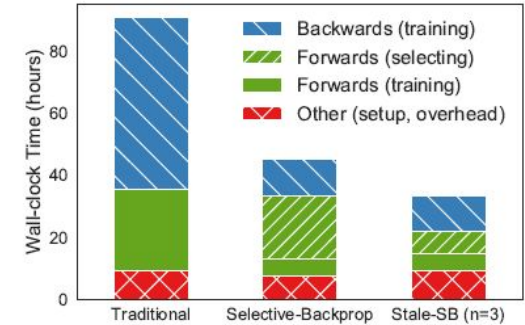


Figure 1: Comparison of selective-backprop approaches [1]

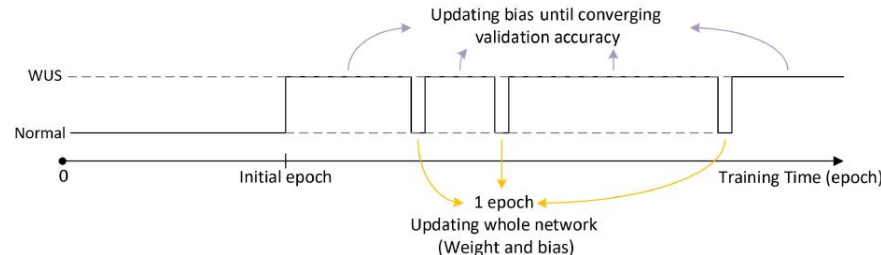


Figure 2: Weight update skipping training cycle [2]

Background Research - Additional Literature

- Adapters [3]
 - Injecting new layers to the network
 - Near identity initialization, trained afterwards
- Movement Pruning [4]
 - Similar to magnitude pruning (prune smallest weight)
 - Both weights with low and high values can be pruned if they shrink in training
- BitFit [5]
 - Similar to weight update skipping
 - No need for normal training step since only doing fine-tuning
 - Update only a small section of biases
- (IA)³ [6]
 - Scales activations by learned vectors

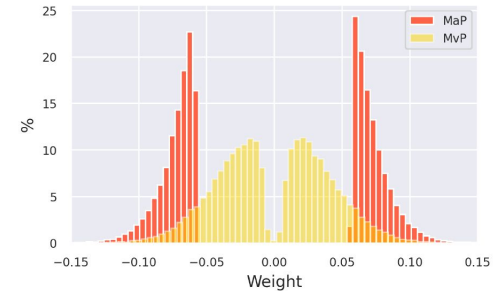


Figure 3: Distribution of the remaining weights after movement pruning [4]

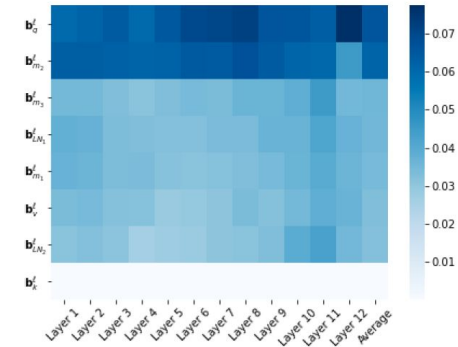


Figure 4: Change in bias components in BitFit [5]

Downstream Task & Dataset

- Sentiment Analysis
 - SST-2 Binary Classification Dataset (GLUE) [7]
- Question Answering
 - SQuAD [8]

Training Loop

- Data preprocessing / tokenization

```

12 def tokenize(dataset, tokenizer, max_length):
13     # If we need to truncate, truncate the context instead of the question
14     tokenized_inputs = tokenizer(
15         dataset["question"],
16         dataset["context"],
17         max_length=max_length,
18         truncation="only_second",
19         padding="max_length",
20         return_offsets_mapping=True
21     )
22
23     offset_mapping = tokenized_inputs.pop("offset_mapping")
24     start_positions = []
25     end_positions = []
26
27     # Update the start and end positions
28     for i, offsets in enumerate(offset_mapping):
29         answer = dataset["answers"][i]
30         # Start/end character index of the answer in the text
31         start_char = answer["answer_start"][0]
32         end_char = start_char + len(answer["text"][0])
33
34         sequence_ids = tokenized_inputs.sequence_ids(i)

```

```

7 def load_dataset(dataset = "squad", split = None):
8     dataset = huggingface_datasets.load_dataset(dataset, split=split)
9     return dataset

```

```

35     # String of 1s in sequence_ids[i] is the context, find first and last
36     context_start = sequence_ids.index(1)
37     context_end = len(sequence_ids) - 1 - sequence_ids[::-1].index(1)
38
39     # If the answer is out of the span (in the question) or after the context, set to 0,0
40     if end_char < offsets[context_start][0] or start_char > offsets[context_end][1]:
41         start_positions.append(0)
42         end_positions.append(0)
43     else:
44         idx = context_start
45
46         while offsets[idx][0] <= start_char and idx < context_end:
47             idx += 1
48         start_positions.append(idx - 1)
49
50         while idx >= context_start and offsets[idx][1] >= end_char:
51             idx -= 1
52         end_positions.append(idx + 1)
53
54     tokenized_inputs["start_positions"] = start_positions
55     tokenized_inputs["end_positions"] = end_positions
56
57     return tokenized_inputs

```

Training Loop

- Layer freezing

```
59 def freeze(model):
60     # Set requires_grad to False for all parameters
61     for param in model.parameters():
62         param.requires_grad = False
63
64
65 def unfreeze_last(model, unfreeze_n = 1):
66     # Unfreeze the last n layers
67     if "distilbert" in model.name_or_path:
68         # DistilBERT
69         for param in model.distilbert.transformer.layer[-unfreeze_n:].parameters():
70             param.requires_grad = True
71     elif "bert" in model.name_or_path:
72         # BERT
73         for param in model.bert.encoder.layer[-unfreeze_n:].parameters():
74             param.requires_grad = True
75
76     # QA head
77     for param in model.qa_outputs.parameters():
78         param.requires_grad = True
```

Training Loop

- Training on the full dataset or a subset

```

81 def train(n_train = None, n_val = None, model_name = "distilbert-base-cased", freeze = True, unfreeze_n = 1):
82     # ===== MODEL ===== #
83     # Load model
84     model = AutoModelForQuestionAnswering.from_pretrained(model_name)
85     # Freeze all but last layer + QA head
86     if freeze:
87         freeze(model)
88         unfreeze_last(model, unfreeze_n=unfreeze_n)
89
90     # ===== DATA ===== #
91     dataset = load_dataset()
92     data_collator = DefaultDataCollator()
93
94     # Take subset of dataset if specified
95     if n_train:
96         dataset["train"] = dataset["train"].select(range(n_train))
97     if n_val:
98         dataset["validation"] = dataset["validation"].select(range(n_val))
99
100     # Tokenize the dataset with our tokenization function
101     tokenizer = AutoTokenizer.from_pretrained(model_name)
102     max_length = model.config.max_position_embeddings
103     tokenize_partial = partial(tokenize, tokenizer=tokenizer, max_length=max_length)
104     tokenized_dataset = dataset.map(tokenize_partial, batched=True, remove_columns=dataset["train"].column_names)
105
106 
```

```

107
108 # ===== TRAINING ===== #
109 training_args = TrainingArguments(
110     output_dir="results",
111     evaluation_strategy="epoch",
112     learning_rate=2e-5,
113     per_device_train_batch_size=16,
114     per_device_eval_batch_size=16,
115     num_train_epochs=5,
116     weight_decay=0.01
117 )
118
119 trainer = Trainer(
120     model,
121     training_args,
122     train_dataset=tokenized_dataset["train"],
123     eval_dataset=tokenized_dataset["validation"],
124     data_collator=data_collator,
125     tokenizer=tokenizer
126 )
127
128 # Perform validation before training
129 print("Evaluating before training (epoch 0)...")
130 metrics = trainer.evaluate()
131 print(metrics)
132
133 trainer.train()
134
135
136 if __name__ == "__main__":
137     train(5000, 500, model_name = "distilbert-base-cased", freeze=False)

```


Generic Update Policy Hook

- UpdatePolicy class
- apply() to enable/disable parameter updates
- Disable parameter by name or freeze entire layers
- Can be extended

```

6 class UpdatePolicy:
7     def __init__(self, model: Module):
8         self.model = model
9
10    best_eval_loss = math.inf
11    unfreeze_n = 1
12
13    # Applies the update policy by setting requires_grad on parameters
14    def apply(self, epoch=0, metrics=None):
15        # Update all parameters
16        for param in self.model.parameters():
17            param.requires_grad = True
18
19        # Disable updates on individual parameter after the first 2 epochs
20        parameter = self.model.get_parameter("distilbert.embeddings.word_embeddings.weight")
21        parameter.requires_grad = epoch in range(2)
22
23        # Freeze all layers but the last one if eval_loss is better than best_eval_loss by a margin
24        if metrics["eval_loss"] < self.best_eval_loss - 1:
25            self.freeze()
26            self.unfreeze_last()
27
28    self.best_eval_loss = min(self.best_eval_loss, metrics["eval_loss"])
29

```

Generic Update Policy Hook

- Manual training loop that injects update policy

```
144 # Perform training
145 for epoch in range(int(training_args.num_train_epochs)):
146     # Apply the update policy before each epoch
147     update_policy.apply(epoch, metrics)
148
149     for batch in train_dataset.iter(training_args.per_device_train_batch_size):
150         batch = {k: torch.tensor(v, dtype=torch.long) for k, v in batch.items()}
151         trainer.training_step(model, batch)
152         trainer.optimizer.step()
153         trainer.optimizer.zero_grad()
154         trainer.lr_scheduler.step()
155         progress_bar.update(1)
156
157     # Evaluate after each epoch
158     print(f"Evaluating epoch {epoch + 1}...")
159     metrics = trainer.evaluate()
160     print(metrics)
```

Open Issues & Solutions

- Which dataset to choose?
 - SQuAD or SST
- Which model to use?
 - DistilBERT?

Plan for the Next Two Weeks

- Decide on our implementation of selective parameter updating
- Start experimentations

References

- [1] A. H. Jiang *et al.*, 'Accelerating Deep Learning by Focusing on the Biggest Losers', *arXiv [cs.LG]*. 2019.
- [2] P. Safayanikoo and I. Akturk, 'Weight Update Skipping: Reducing Training Time for Artificial Neural Networks', *arXiv [cs.LG]*. 2020.
- [3] N. Houlsby *et al.*, 'Parameter-Efficient Transfer Learning for NLP', *arXiv [cs.LG]*. 2019.
- [4] V. Sanh, T. Wolf, and A. M. Rush, 'Movement Pruning: Adaptive Sparsity by Fine-Tuning', *arXiv [cs.CL]*. 2020.
- [5] E. B. Zaken, S. Ravfogel, and Y. Goldberg, 'BitFit: Simple Parameter-efficient Fine-tuning for Transformer-based Masked Language-models', *arXiv [cs.LG]*. 2022.
- [6] H. Liu *et al.*, 'Few-Shot Parameter-Efficient Fine-Tuning is Better and Cheaper than In-Context Learning', *arXiv [cs.LG]*. 2022.
- [7] A. Wang, A. Singh, J. Michael, F. Hill, O. Levy, and S. R. Bowman, 'GLUE: A Multi-Task Benchmark and Analysis Platform for Natural Language Understanding', *arXiv [cs.CL]*. 2019.
- [8] P. Rajpurkar, J. Zhang, K. Lopyrev, and P. Liang, 'SQuAD: 100,000+ Questions for Machine Comprehension of Text', *arXiv [cs.CL]*. 2016.