

Analisi del Fine Tuning del Modello BERT per l'Analisi di Risposte LLM, Versione 1.8

10 aprile 2025

Sommario

Questo documento presenta un'analisi dettagliata del processo di fine tuning effettuato su un modello BERT (`bert-base-uncased`) per l'analisi di risposte generate da modelli LLM. Vengono descritti il dataset utilizzato, la metodologia di training, l'andamento della loss, ed i risultati quantitativi ottenuti in fase di valutazione. In coda al documento è incluso il codice completo di training.

1 Introduzione

Il fine tuning di modelli pre-addestrati come BERT permette di adattare la rappresentazione linguistica a task specifici, come in questo caso l'analisi delle risposte generate da modelli LLM. L'approccio adottato prevede l'uso di un dataset suddiviso in set di train, validazione e test. Il modello è stato inizializzato con i pesi di `bert-base-uncased` e successivamente sottoposto a ulteriori allenamenti sul dataset target, aggiornando anche il classificatore (con parametri `classifier.weight` e `classifier.bias`) che non erano presenti nel checkpoint di partenza.

2 Dataset e Pre-Processing

Il dataset utilizzato è organizzato in un `DatasetDict` con le seguenti caratteristiche:

- **Training:** 12.000 esempi (poi suddivisi in 9.600 per l'addestramento e 2.400 per la validazione)
- **Test:** 1.773 esempi

Le colonne principali del dataset sono:

- **response:** Contiene le risposte generate dai modelli LLM.
- **label:** Indicatore della classe (con etichette uniche: "0" e "1").

Il pre-processing ha incluso la mappatura delle etichette sui rispettivi ID, la tokenizzazione dei testi con un tokenizer pre-addestrato e la rimozione di colonne non utili all'addestramento.

3 Dettagli del Training e Risultati

3.1 Impostazioni e Log del Training

Il training è avvenuto per 3 epoche con i seguenti parametri chiave:

- **Learning rate:** 2×10^{-5}
- **Batch size:** 16 per dispositivo (sia per train che per eval)
- **Warmup steps:** 500
- **Weight decay:** 0.01
- **FP16:** Abilitato per l'ottimizzazione
- **Gradient Accumulation:** 2 steps

Il training log mostra un andamento della **loss** decrescente. Alcuni step rappresentativi sono:

- **Step 100:** Loss ≈ 0.6206
- **Step 300:** Loss ≈ 0.2167
- **Step 700:** Loss ≈ 0.0881
- **Step 900:** Loss ≈ 0.0752

Questi valori indicano una convergenza adeguata del modello nel corso del training, con una progressiva riduzione dell'errore.

3.2 Valutazione sul Test Set

I risultati finali sul test set evidenziano un'alta capacità predittiva del modello:

- **Loss:** 0.501045
- **Accuracy:** 90.02%
- **Precision:** 90.34%
- **Recall:** 90.02%
- **F1-score:** 89.92%

Questi risultati suggeriscono che il modello, dopo il fine tuning, sia in grado di catturare efficacemente le caratteristiche del problema di classificazione. L'accuratezza elevata e i valori bilanciati di precision e recall indicano una buona generalizzazione, sebbene una loss non trascurabile possa ancora essere oggetto di ulteriori ottimizzazioni, per esempio tramite tecniche di regolarizzazione o un ulteriore tuning degli iperparametri.

3.3 Considerazioni Tecniche

- Il messaggio "Some weights of BertForSequenceClassification were not initialized..." indica che il classificatore è stato inizializzato casualmente e quindi ha dovuto apprendere da zero la parte finale della rete.
- L'approccio di suddividere il dataset in training, validazione e test consente una stima affidabile della generalizzazione del modello.
- L'uso di metriche multiple (accuracy, precision, recall e F1-score) offre una visione completa delle performance, importante in applicazioni di NLP dove la distribuzione delle classi può essere sbilanciata.

4 Conclusioni

Il fine tuning effettuato sul modello BERT ha prodotto un sistema capace di classificare con una buona accuratezza le risposte generate dai modelli LLM. I risultati quantitativi confermano l'efficacia dell'approccio adottato e la corretta impostazione dei parametri di training. Per futuri miglioramenti si potrebbero esplorare:

- Ulteriore ottimizzazione degli iperparametri;
- Tecniche di data augmentation per incrementare la robustezza del modello;
- Strategie di regolarizzazione per ridurre ulteriormente la loss.

Codice di Training

Di seguito viene riportato il codice utilizzato per il fine tuning del modello:

```
1 ! pip install transformers datasets evaluate torch accelerate -
  U
2 !pip install -U transformers
3
4 # 'accelerate' è raccomandato per Trainer per ottimizzare l'uso
  della GPU/TPU
5 import os
```

```

6 import numpy as np
7 import pandas as pd
8 from datasets import load_dataset
9 from transformers import AutoTokenizer,
    AutoModelForSequenceClassification, TrainingArguments,
    Trainer, TrainerCallback
10 import evaluate
11 from sklearn.metrics import precision_recall_fscore_support
12
13 # Disabilitiamo wandb in modalita offline (salva i log
    localmente)
14 os.environ["WANDB_MODE"] = "offline"
15
16 # Definiamo la directory di output dove salvare il modello e il
    tokenizer
17 output_directory = "my-bert-fine-tuned-model"
18
19 # Carichiamo il dataset
20 data_files = {
21     "train": "dataset_completo.json",
22     "test": "Test2.json"
23 }
24 dataset = load_dataset('json', data_files=data_files)
25
26 print("Struttura del dataset:")
27 print(dataset)
28 print("\nColonne nel dataset train:")
29 print(dataset["train"].column_names)
30 print("\nPrimo esempio nel dataset:")
31 print(dataset["train"][0])
32
33 # Determiniamo le colonne di testo e etichette
34 first_example = dataset["train"][0]
35 text_column = None
36 label_column = None
37
38 # Trova la colonna del testo (quella più lunga)
39 longest_text_len = 0
40 for col in first_example:
41     if isinstance(first_example[col], str) and len(
        first_example[col]) > longest_text_len:
42         longest_text_len = len(first_example[col])
43         text_column = col
44
45 # Trova la colonna delle etichette (cerca 'label', 'class' o '
    category')
46 for col in first_example:
47     if 'label' in col.lower() or 'class' in col.lower() or '
        category' in col.lower():
48         label_column = col
49         break
50
51 if text_column is None:
52     raise ValueError("Non è stata trovata una colonna di testo.")

```

```

        Specifica manualmente il nome della colonna.")
53 if label_column is None:
54     # Se non troviamo una colonna di etichette evidente,
        utilizziamo una colonna non di testo
55     for col in first_example:
56         if col != text_column and not isinstance(first_example[
            col], str):
57             label_column = col
58             break
59
60 print(f"\nColonna di testo identificata: '{text_column}'")
61 print(f"Colonna di etichette identificata: '{label_column}'")
62
63 # Pre-processiamo le etichette e creiamo il mapping label -> ID
64 def get_unique_labels(examples):
65     labels = examples[label_column]
66     unique_labels = set()
67     for label in labels:
68         if isinstance(label, list):
69             for l in label:
70                 unique_labels.add(l)
71         else:
72             unique_labels.add(label)
73     return list(unique_labels)
74
75 unique_labels = get_unique_labels(dataset["train"])
76 print(f"\nEtichette uniche trovate: {unique_labels}")
77
78 label_to_id = {label: i for i, label in enumerate(sorted(
    unique_labels))}
79 id_to_label = {i: label for label, i in label_to_id.items()}
80 print(f"\nMapping etichette -> ID: {label_to_id}")
81
82 def preprocess_labels(examples):
83     result = dict(examples)
84     labels = examples[label_column]
85     processed_labels = []
86     for label in labels:
87         if isinstance(label, list):
88             processed_labels.append(label_to_id[label[0]] if
                label else 0)
89         else:
90             processed_labels.append(label_to_id[label])
91     result[label_column] = processed_labels
92     return result
93
94 processed_dataset = dataset.map(preprocess_labels, batched=True
    )
95
96 # Scegliamo il modello e il tokenizer (ad es. "bert-base-
    uncased")
97 model_checkpoint = "bert-base-uncased"
98 tokenizer = AutoTokenizer.from_pretrained(model_checkpoint)
99

```

```

100 def tokenize_function(examples):
101     return tokenizer(
102         examples[text_column],
103         padding="max_length",
104         truncation=True,
105         max_length=128
106     )
107
108 # Tokenizziamo il dataset
109 tokenized_datasets = processed_dataset.map(tokenize_function,
110     batched=True)
111 tokenized_datasets = tokenized_datasets.remove_columns([col for
112     col in processed_dataset["train"].column_names if col !=
113     label_column])
114 tokenized_datasets = tokenized_datasets.rename_column(
115     label_column, "labels")
116 tokenized_datasets.set_format("torch")
117
118 # Creiamo i set di training, validazione e test
119 train_testvalid = tokenized_datasets["train"].train_test_split(
120     test_size=0.2, seed=42)
121 train_dataset = train_testvalid["train"]
122 validation_dataset = train_testvalid["test"]
123 test_dataset = tokenized_datasets["test"]
124
125 print(f"\nDimensione del dataset di training completo: {len(
126     train_dataset)} esempi")
127 print(f"Dimensione del dataset di validazione: {len(
128     validation_dataset)} esempi")
129 print(f"Dimensione del dataset di test: {len(test_dataset)}
130     esempi")
131
132 # Impostiamo la metrica di accuracy
133 metric = evaluate.load("accuracy")
134
135 def compute_metrics(eval_pred):
136     logits, labels = eval_pred
137     predictions = np.argmax(logits, axis=-1)
138     accuracy = metric.compute(predictions=predictions,
139         references=labels)
140     precision, recall, f1, _ = precision_recall_fscore_support(
141         labels, predictions, average='weighted')
142     return {
143         'accuracy': accuracy['accuracy'],
144         'precision': precision,
145         'recall': recall,
146         'f1': f1
147     }
148
149 num_labels = len(label_to_id)
150 model = AutoModelForSequenceClassification.from_pretrained(
151     model_checkpoint,
152     num_labels=num_labels
153 )

```

```

144
145 # Definiamo un callback personalizzato per salvare e stampare i
      log di training
146 class LogCallback(TrainerCallback):
147     def __init__(self):
148         self.logs = [] # Lista per salvare i log intermedi
149
150     def on_log(self, args, state, control, logs=None, **kwargs):
151         :
152         if logs is not None:
153             # Salviamo solo i log rilevanti (es. loss, lr, step
              , epoch)
154             self.logs.append({
155                 'step': state.global_step,
156                 'epoch': state.epoch,
157                 'loss': logs.get('loss', None),
158                 'learning_rate': logs.get('learning_rate', None)
159             },
160             'eval_loss': logs.get('eval_loss', None)
161             })
162
163     def on_train_end(self, args, state, control, **kwargs):
164         # Alla fine dell'addestramento stampiamo una tabella
          riassuntiva
165         df = pd.DataFrame(self.logs)
166         print("\n=== Riassunto Training Log ===")
167         # Stampiamo log ogni 100 step
168         df_summary = df[df['step'] % 100 == 0]
169         print(df_summary.to_string(index=False))
170
171 # Configuriamo gli argomenti di addestramento
172 training_args = TrainingArguments(
173     output_dir=output_directory,
174     eval_steps=100, # Valutazione ogni 100 step
175     save_steps=100, # Salvataggio ogni 100 step
176     logging_steps=100, # Stampa dei log ogni 100
              step
177     learning_rate=2e-5,
178     per_device_train_batch_size=16,
179     per_device_eval_batch_size=16,
180     num_train_epochs=3,
181     weight_decay=0.01,
182     warmup_steps=500,
183     fp16=True,
184     gradient_accumulation_steps=2,
185     save_total_limit=2,
186     report_to="none"
187 )
188
189 # Inizializziamo il trainer passando il parametro aggiornato "
      processing_class" invece di "tokenizer"
190 log_callback = LogCallback()
191 trainer = Trainer(
192     model=model,

```

```

191     args=training_args,
192     train_dataset=train_dataset,
193     eval_dataset=validation_dataset,
194     processing_class=tokenizer, # Utilizziamo il nuovo
        parametro in sostituzione di 'tokenizer'
195     compute_metrics=compute_metrics,
196     callbacks=[log_callback]
197 )
198
199 print("Inizio addestramento sull'intero dataset...")
200 trainer.train()
201 print("Addestramento completato!")
202
203 # Valutazione sul test set
204 print("Valutazione sul test set completo...")
205 test_results = trainer.evaluate(test_dataset)
206 print("Risultati test:", test_results)
207
208 # Salva il modello e il tokenizer
209 trainer.save_model(output_directory)
210 tokenizer.save_pretrained(output_directory)
211 print(f"Modello e tokenizer salvati in {output_directory}")
212
213 # Stampa il dizionario delle etichette per uso futuro
214 print("\nDizionario delle etichette (utile per interpretare le
        previsioni):")
215 print(id_to_label)
216
217 # Esempio di utilizzo del modello:
218 print("\nEsempio di utilizzo del modello:")
219 print('from transformers import
        AutoModelForSequenceClassification, AutoTokenizer')
220 print(f'model = AutoModelForSequenceClassification.
        from_pretrained("{output_directory}")')
221 print(f'tokenizer = AutoTokenizer.from_pretrained("{
        output_directory}")')
222 print('inputs = tokenizer("Esempio di testo", return_tensors="
        pt")')
223 print('outputs = model(**inputs)')
224 print('predictions = outputs.logits.argmax(-1).item()')
225 print('etichetta_prevista = id_to_label[predictions] #
        Converti l\'ID nell\'etichetta originale')
226
227 # Stampiamo una tabella finale con i risultati complessivi
228 summary_dict = {
229     "Metric": ["eval_loss", "accuracy", "precision", "recall",
        "f1"],
230     "Valore": [
231         test_results.get('eval_loss', 'N/A'),
232         test_results.get('eval_accuracy', 'N/A'),
233         test_results.get('eval_precision', 'N/A'),
234         test_results.get('eval_recall', 'N/A'),
235         test_results.get('eval_f1', 'N/A')
236 ]

```



```
237 }  
238 df_summary = pd.DataFrame(summary_dict)  
239 print("\n== Tabella Riassuntiva dei Risultati dell' "  
    Adestramento ==")  
240 print(df_summary.to_string(index=False))
```

Listing 1: Codice di Training per il Fine Tuning di BERT