

TaskB_DistilBERT

November 27, 2024

```
[11]: import evaluate
import transformers
import numpy as np
import pandas as pd
from datasets import load_dataset, Dataset, DatasetDict
from transformers import (
    AutoTokenizer,
    AutoModelForSequenceClassification,
    TrainingArguments,
    Trainer,
    DistilBertForSequenceClassification,
)
```

```
[12]: """
Download dataset SubtaskA.jsonl from
https://github.com/mbzuai-nlp/M4GT-Bench.
"""

DATA_PATH = "C:/Users/Admin/Desktop/cse847_proj/SubtaskB.jsonl"

# initialize dataset
df = pd.read_json(DATA_PATH, lines=True)
df = df[['text', 'label', 'model']]
dataset = Dataset.from_pandas(df)

# split dataset
a = dataset.train_test_split(test_size=0.20)
b = a['test'].train_test_split(test_size=0.5)
dataset = DatasetDict({
    'train': a['train'],
    'valid': b['train'],
    'test': b['test'],
})
print(dataset)
```

```
DatasetDict({
  train: Dataset({
    features: ['text', 'label', 'model'],
    num_rows: 98248
```

```

    })
    valid: Dataset({
        features: ['text', 'label', 'model'],
        num_rows: 12281
    })
    test: Dataset({
        features: ['text', 'label', 'model'],
        num_rows: 12282
    })
})

```

```
[13]: print(df.label.unique())
      df[['text', 'label']]
```

```
[2 1 0 3 5 4 6]
```

```
[13]:
```

	text	label
0	We consider a system of many polymers in solut...	2
1	We present a catalog of 66 YSOs in the Serpens...	2
2	Spectroscopic Observations of the Intermediate...	2
3	We present a new class of stochastic Lie group...	2
4	ALMA as the ideal probe of the solar chromosph...	2
...
122806	Title: The Unsung Heroes: Seagoing Cowboys and...	0
122807	Title: The Benefits of Autonomy: Student-led P...	0
122808	The Electoral College system, established by t...	0
122809	In the ever-evolving landscape of education, c...	0
122810	When faced with critical decisions, the wise o...	0

```
[122811 rows x 2 columns]
```

```
[14]: """
      Initialize tokenizer and model.
      """
      model_id = "distilbert-base-uncased"

      # init tokenizer
      tokenizer = AutoTokenizer.from_pretrained(model_id)

      # init model
      model = DistilBertForSequenceClassification.from_pretrained(
          model_id,
          num_labels=7,
      )

```

Some weights of DistilBertForSequenceClassification were not initialized from the model checkpoint at distilbert-base-uncased and are newly initialized: ['classifier.bias', 'classifier.weight', 'pre_classifier.bias', 'pre_classifier.weight']

You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.

```
[15]: """
      Tokenize dataset.
      """
      def tokenize(X):
          return tokenizer(
              X["text"],
              padding="max_length",
              truncation=True,
              return_tensors="pt",
          )

      # tokenize data
      tokenized_datasets = dataset.map(tokenize, batched=True)
      print(tokenized_datasets)
```

```
Map:   0%|          | 0/98248 [00:00<?, ? examples/s]
Map:   0%|          | 0/12281 [00:00<?, ? examples/s]
Map:   0%|          | 0/12282 [00:00<?, ? examples/s]

DatasetDict({
  train: Dataset({
    features: ['text', 'label', 'model', 'input_ids', 'attention_mask'],
    num_rows: 98248
  })
  valid: Dataset({
    features: ['text', 'label', 'model', 'input_ids', 'attention_mask'],
    num_rows: 12281
  })
  test: Dataset({
    features: ['text', 'label', 'model', 'input_ids', 'attention_mask'],
    num_rows: 12282
  })
})
```

```
[16]: """
      Create dataset splits.
      """
      seed = 777
      n_samples = 10_000
      n_test = 1000

      train_dataset = tokenized_datasets["train"].shuffle(seed=seed).
      ↪select(range(n_samples))
```

```
valid_dataset = tokenized_datasets["valid"].shuffle(seed=seed).
    ↪select(range(n_test))
test_dataset = tokenized_datasets["test"].shuffle(seed=seed).
    ↪select(range(n_test))
```

```
[17]: """
      Create Trainer.
      """
      # define metric
      metric = evaluate.load("accuracy")

      def compute_metrics(eval_pred):
          logits, labels = eval_pred
          predictions = np.argmax(logits, axis=-1)
          return metric.compute(predictions=predictions, references=labels)

      # training args
      training_args = TrainingArguments(
          output_dir="C:/Users/Admin/Desktop/cse847_proj/",
          eval_strategy="epoch",
          save_total_limit=2,
      )

      # init trainer
      trainer = Trainer(
          model=model,
          args=training_args,
          train_dataset=train_dataset,
          eval_dataset=valid_dataset,
          compute_metrics=compute_metrics,
      )
```

```
[18]: """
      Train model.
      """
      trainer.train()
```

<IPython.core.display.HTML object>

```
[18]: TrainOutput(global_step=3750, training_loss=0.4712401387532552,
metrics={'train_runtime': 32555.0664, 'train_samples_per_second': 0.922,
'train_steps_per_second': 0.115, 'total_flos': 3974376314880000.0, 'train_loss':
0.4712401387532552, 'epoch': 3.0})
```

```
[19]: """
      Evaluate trained model.
      """
      trainer.evaluate(test_dataset)
```

<IPython.core.display.HTML object>

```
[19]: {'eval_loss': 0.5755491256713867,  
      'eval_accuracy': 0.873,  
      'eval_runtime': 299.7426,  
      'eval_samples_per_second': 3.336,  
      'eval_steps_per_second': 0.417,  
      'epoch': 3.0}
```

```
[20]: """  
      Summarize model.  
      """  
      print(model)
```

```
DistilBertForSequenceClassification(  
  (distilbert): DistilBertModel(  
    (embeddings): Embeddings(  
      (word_embeddings): Embedding(30522, 768, padding_idx=0)  
      (position_embeddings): Embedding(512, 768)  
      (LayerNorm): LayerNorm((768,), eps=1e-12, elementwise_affine=True)  
      (dropout): Dropout(p=0.1, inplace=False)  
    )  
    (transformer): Transformer(  
      (layer): ModuleList(  
        (0-5): 6 x TransformerBlock(  
          (attention): DistilBertSdpaAttention(  
            (dropout): Dropout(p=0.1, inplace=False)  
            (q_lin): Linear(in_features=768, out_features=768, bias=True)  
            (k_lin): Linear(in_features=768, out_features=768, bias=True)  
            (v_lin): Linear(in_features=768, out_features=768, bias=True)  
            (out_lin): Linear(in_features=768, out_features=768, bias=True)  
          )  
          (sa_layer_norm): LayerNorm((768,), eps=1e-12, elementwise_affine=True)  
          (ffn): FFN(  
            (dropout): Dropout(p=0.1, inplace=False)  
            (lin1): Linear(in_features=768, out_features=3072, bias=True)  
            (lin2): Linear(in_features=3072, out_features=768, bias=True)  
            (activation): GELUActivation()  
          )  
          (output_layer_norm): LayerNorm((768,), eps=1e-12,  
elementwise_affine=True)  
        )  
      )  
    )  
    (pre_classifier): Linear(in_features=768, out_features=768, bias=True)  
    (classifier): Linear(in_features=768, out_features=7, bias=True)  
    (dropout): Dropout(p=0.2, inplace=False)
```

)

[]: