

Homework 3

- Some questions require writing Python code and computing results, and the rest of them have written answers. For coding problems, you will have to fill out all code blocks that say **YOUR CODE HERE**.
- You will need to use GPU, which can be added through `Edit > Notebook Settings > Hardware accelerator > (GPU)`
- For text-based answers, you should replace the text that says **Write your answer here...** with your actual answer.
- This assignment is designed such that each cell takes a few minutes (if that) to run. If it is taking longer than that, you might have made a mistake in your code.

How to submit this problem set:

- Write all the answers in this Colab notebook. Once you are finished, generate a PDF via (`File -> Print -> Save as PDF`) and upload it to Gradescope.
- **Important:** check your PDF before you submit to Gradescope to make sure it exported correctly. If Colab gets confused about your syntax, it will sometimes terminate the PDF creation routine early.
- When creating your final version of the PDF to hand in, please do a fresh restart and execute every cell in order. One handy way to do this is by clicking `Runtime -> Run All` in the notebook menu.

<https://colab.research.google.com/drive/1tTcRE6fvN6mQsekDMYM-v9hL4nP8QLIs?usp=sharing>

▼ Part 0: Setup

▼ Installing Hugging Face's Transformers and Additional Libraries

We will use Hugging Face's Transformers (<https://github.com/huggingface/transformers>), an open-source library that provides general-purpose architectures for natural language understanding and generation with a collection of various pretrained models made by the NLP community. This library will allow us to easily use pretrained models like BERT and perform experiments on top of them. We can use these models to solve downstream target tasks, such as text classification, question answering, sequence labeling, and text generation.

Run the following cell to install Hugging Face's Transformers library and some other useful tools. This cell will also download data used later in the assignment.

```

1 ! nvidia-smi
2 !pip install -q transformers==4.17.0 datasets==2.0.0 rich[jupyter]
3 !pip install -q googletrans==3.1.0a0
4 !pip install -q -U PyDrive
5 !apt install jq
6
7 import torch
8 from pydrive.auth import GoogleAuth
9 from pydrive.drive import GoogleDrive
10 from google.colab import auth
11 from oauth2client.client import GoogleCredentials
12
13 import os
14 import zipfile
15 import collections
16
17 import pandas as pd
18 import numpy as np
19 from matplotlib import pyplot as plt
20 from matplotlib import cm
21 from rich import print as rich_print
22
23 from transformers import GPT2LMHeadModel, GPT2Tokenizer

```

Processes:						
GPU	GI	CI	PID	Type	Process name	GPU Memory Usage
ID	ID	ID				
=====						
No running processes found						

```

| 3.8 MB 35.0 MB/s
| 325 kB 77.6 MB/s
| 237 kB 76.8 MB/s
| 182 kB 75.1 MB/s
|
| 7.6 MB 68.9 MB/s
| 880 kB 62.6 MB/s
| 212 kB 80.4 MB/s
| 132 kB 68.5 MB/s
| 127 kB 66.2 MB/s
| 51 kB 8.4 MB/s
| 1.6 MB 57.9 MB/s
Building wheel for sacremoses (setup.py) ... done
| 55 kB 3.6 MB/s
| 1.5 MB 66.4 MB/s
| 42 kB 1.3 MB/s
| 65 kB 4.5 MB/s
| 52 kB 37.1 MB/s

```

```

Building wheel for googlettrans (setup.py) ... done
Reading package lists... Done
Building dependency tree
Reading state information... Done
The following package was automatically installed and is no longer required:
  libnvidia-common-460
Use 'apt autoremove' to remove it.
The following additional packages will be installed:
  libjq1 libonig4
The following NEW packages will be installed:
  jq libjq1 libonig4
0 upgraded, 3 newly installed, 0 to remove and 7 not upgraded.
Need to get 276 kB of archives.
After this operation, 930 kB of additional disk space will be used.
Get:1 http://archive.ubuntu.com/ubuntu bionic/universe amd64 libonig4 amd64 6.7.0-1 [
Get:2 http://archive.ubuntu.com/ubuntu bionic/universe amd64 libjq1 amd64 1.5+dfsg-2
Get:3 http://archive.ubuntu.com/ubuntu bionic/universe amd64 jq amd64 1.5+dfsg-2 [45.
Fetched 276 kB in 2s (166 kB/s)
Selecting previously unselected package libonig4:amd64.
(Reading database ... 124015 files and directories currently installed.)
Preparing to unpack .../libonig4_6.7.0-1_amd64.deb ...
Unpacking libonig4:amd64 (6.7.0-1) ...
Selecting previously unselected package libjq1:amd64.
Preparing to unpack .../libjq1_1.5+dfsg-2_amd64.deb ...
Unpacking libjq1:amd64 (1.5+dfsg-2) ...
Selecting previously unselected package jq.
Preparing to unpack .../jq_1.5+dfsg-2_amd64.deb ...
Unpacking jq (1.5+dfsg-2) ...
Setting up libonig4:amd64 (6.7.0-1) ...
Setting up libjq1:amd64 (1.5+dfsg-2) ...
Setting up jq (1.5+dfsg-2) ...
Processing triggers for man-db (2.8.3-2ubuntu0.1) ...

```

```

1 def tsv_to_csv(in_file, out_file):
2     data = pd.read_csv(in_file, sep='\t')
3     data.to_csv(out_file, sep=',', index=False)

1 auth.authenticate_user()
2 gauth = GoogleAuth()
3 gauth.credentials = GoogleCredentials.get_application_default()
4 drive = GoogleDrive(gauth)
5 print('success!')
6
7 data_file = drive.CreateFile({'id': '1zeo8FcaNUnhN660mGMNEAPvxOE4DPOnE'})
8 data_file.GetContentFile('hw1.zip')
9
10 with zipfile.ZipFile('hw1.zip', 'r') as zip_file:
11     zip_file.extractall('./')
12 os.remove('hw1.zip')
13 os.chdir('hw1')
14 print("Data and supporting code downloaded!")

```

```

15
16 tsv_to_csv('data/tinySST/dev.tsv', 'data/tinySST/dev.csv')
17 tsv_to_csv('data/tinySST/train.tsv', 'data/tinySST/train.csv')
18 print('finished preprocessing data')
19
20 pretrained_models_dir = './pretrained_models_dir'
21 if not os.path.isdir(pretrained_models_dir):
22     os.mkdir(pretrained_models_dir)
23 print('model directory created')
24
25 !pip install -q -r requirements.txt
26 print('extra packages installed!')
27
28 !wget -nv https://nlp.stanford.edu/data/coqa/coqa-train-v1.0.json
29 !wget -nv https://nlp.stanford.edu/data/coqa/coqa-dev-v1.0.json
30 !cat coqa-dev-v1.0.json | jq '{data: [.data[] | del(.answers[].bad_turn) | del(.question
31 !cat coqa-train-v1.0.json | jq '{data: [.data[] | del(.answers[].bad_turn) | del(.questi
32 print('Download coqa dataset!')

```

success!
Data and supporting code downloaded!
finished preprocessing data
model directory created

	43 kB	1.6 MB/s
	981 kB	31.5 MB/s

Building wheel for seqeval (setup.py) ... done
Building wheel for langdetect (setup.py) ... done
extra packages installed!
2022-12-03 15:08:01 URL: <https://downloads.cs.stanford.edu/nlp/data/coqa/coqa-train-v1.0>
2022-12-03 15:08:05 URL: <https://downloads.cs.stanford.edu/nlp/data/coqa/coqa-dev-v1.0.j>
Download coqa dataset!

```

1 tokenizer = GPT2Tokenizer.from_pretrained("gpt2")
2 model = GPT2LMHeadModel.from_pretrained("gpt2", pad_token_id=tokenizer.eos_token_id)

```

Downloading: 100%	0.99M/0.99M [00:01<00:00, 953kB/s]
Downloading: 100%	446k/446k [00:01<00:00, 450kB/s]
Downloading: 100%	665/665 [00:00<00:00, 18.9kB/s]
Downloading: 100%	523M/523M [00:13<00:00, 34.9MB/s]

▼ Part 1. Beam Search

We're going to explore decoding from a pretrained GPT-2 model using beam search. Run the below cell to set up some beam search utilities.

```

1 def init_beam_search(model, input_ids, num_beams):

```

```

2     assert len(input_ids.shape) == 2
3     beam_scores = torch.zeros(num_beams, dtype=torch.float32, device=model.device)
4     beam_scores[1:] = -1e9 # Break ties in first round.
5     new_input_ids = input_ids.repeat_interleave(num_beams, dim=0).to(model.device)
6     return new_input_ids, beam_scores
7
8 def run_beam_search(model, tokenizer, input_text, num_beams=5, num_decode_steps=10, sco
9
10    input_ids = tokenizer.encode(input_text, return_tensors='pt')
11    input_ids, beam_scores = init_beam_search(model, input_ids, num_beams)
12    token_scores = beam_scores.clone().view(num_beams, 1)
13
14    model_kwargs = {}
15    for i in range(num_decode_steps):
16        model_inputs = model.prepare_inputs_for_generation(input_ids, **model_kwargs) #
17        outputs = model(**model_inputs, return_dict=True)
18        next_token_logits = outputs.logits[:, -1, :]
19        vocab_size = next_token_logits.shape[-1]
20        this_token_scores = torch.log_softmax(next_token_logits, -1)
21
22        # Process token scores.
23        processed_token_scores = this_token_scores
24        for processor in score_processors:
25            processed_token_scores = processor(input_ids, processed_token_scores)
26
27        # Update beam scores.
28        next_token_scores = processed_token_scores + beam_scores.unsqueeze(-1)
29        next_token_scores = next_token_scores.view(num_beams * vocab_size)
30
31        # Find top-scoring beams.
32        next_token_scores, next_tokens = torch.topk(
33            next_token_scores, num_beams, dim=0, largest=True, sorted=True
34        )
35
36        # Transform tokens since we reshaped earlier.
37        next_indices = torch.div(next_tokens, vocab_size, rounding_mode="floor") # This
38        next_tokens = next_tokens % vocab_size
39
40        # Update tokens and beam scores.
41        input_ids = torch.cat([input_ids[next_indices, :], next_tokens.unsqueeze(-1)], d
42        beam_scores = next_token_scores
43
44        # UNCOMMENT: To use original scores instead.
45        # token_scores = torch.cat([token_scores[next_indices, :], this_token_scores[nex
46        token_scores = torch.cat([token_scores[next_indices, :], processed_token_scores[
47
48        # Update hidden state.
49        model_kwargs = model._update_model_kwargs_for_generation(outputs, model_kwargs,
50        model_kwargs["past"] = model._reorder_cache(model_kwargs["past"], next_indices)
51
52    def transfer(x):

```

```

53     return x.cpu() if to_cpu else x
54
55     return {
56         "output_ids": transfer(input_ids),
57         "beam_scores": transfer(beam_scores),
58         "token_scores": transfer(token_scores)
59     }
60
61
62 def run_beam_search(*args, **kwargs):
63     with torch.inference_mode():
64         return run_beam_search(*args, **kwargs)

```



```

1 # Add support for colored printing and plotting.
2 RICH_x = np.linspace(0.0, 1.0, 50)
3 RICH_rgb = (cm.get_cmap(plt.get_cmap('RdYlBu'))(RICH_x)[: , :3] * 255).astype(np.int32)[r
4
5 def print_with_probs(words, probs, prefix=None):
6     def fmt(x, p, is_first=False):
7         ix = int(p * RICH_rgb.shape[0])
8         r, g, b = RICH_rgb[ix]
9         if is_first:
10             return f'[bold rgb(0,0,0) on rgb({r},{g},{b})]{x}'
11         else:
12             return f'[bold rgb(0,0,0) on rgb({r},{g},{b})] {x}'
13     output = []
14     if prefix is not None:
15         output.append(prefix)
16     for i, (x, p) in enumerate(zip(words, probs)):
17         output.append(fmt(x, p, is_first=i == 0))
18     rich_print(''.join(output))

```

▼ Question 1.1 (5 points)

Run the cell below. It produces a sequence of tokens using beam search and the provided prefix.

```

1 num_beams = 5
2 num_decode_steps = 10
3 input_text = 'The brown fox jumps'
4
5 beam_output = run_beam_search(model, tokenizer, input_text, num_beams=num_beams, num_dec
6 for i, tokens in enumerate(beam_output['output_ids']):
7     score = beam_output['beam_scores'][i]
8     print(i, round(score.item() / tokens.shape[-1], 3), tokenizer.decode(tokens, skip_sp

```



```

0 -1.106 The brown fox jumps out of the fox's mouth, and the fox
1 -1.168 The brown fox jumps out of the fox's cage, and the fox
2 -1.182 The brown fox jumps out of the fox's mouth and starts to run

```

```
3 -1.192 The brown fox jumps out of the fox's mouth and begins to lick
4 1 100 The brown fox jumps out of the fox's mouth and begins to lick
```

To get you more acquainted with the code, let's do a simple exercise first. Write your own code in the cell below to generate 3 tokens with a beam size of 4, and then print out the **third most probable** output sequence found during the search. Use the same prefix as above.

```
1 input_text = 'The brown fox jumps'
2
3 # WRITE YOUR CODE HERE!
4 beam_output = run_beam_search(model, tokenizer, input_text, num_beams=4, num_decode_step
5 beam_output
6 i = 2 # 0 1 2 3
7 score = beam_output['beam_scores'][i]
8 tokens = beam_output['output_ids'][i]
9 print(i, round(score.item() / tokens.shape[-1], 3), tokenizer.decode(tokens, skip_specia

2 -0.627 The brown fox jumps up and down
```

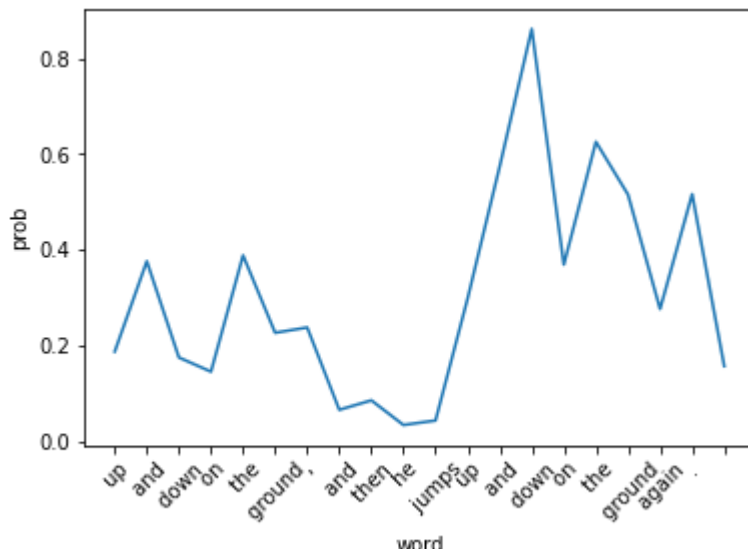
▼ Question 1.2 (10 points)

Run the cells below to visualize the probabilities the model assigns for each generated word when using beam search with beam size 1 (i.e., greedy decoding).

```
1 input_text = 'The brown fox jumps'
2 beam_output = run_beam_search(model, tokenizer, input_text, num_beams=1, num_decode_step
3 probs = beam_output['token_scores'][0, 1:].exp()
4 output_subwords = [tokenizer.decode(tok, skip_special_tokens=True) for tok in beam_output

1 print('Visualization with plot:')
2
3 fig, ax = plt.subplots()
4 plt.plot(range(len(probs)), probs)
5 ax.set_xticks(range(len(probs)))
6 ax.set_xticklabels(output_subwords[-len(probs):], rotation = 45)
7 plt.xlabel('word')
8 plt.ylabel('prob')
9 plt.show()
10
11 print('Visualization with colored text (red for lower probability, and blue for higher):
12
13 print_with_probs(output_subwords[-len(probs):], probs, ' '.join(output_subwords[:-len(pr
```

Visualization with plot:



Why does the model assign a higher probability to tokens generated later than to tokens generated earlier?

WRITE YOUR ANSWER HERE IN A FEW SENTENCES

At first, the model only has 4 words "The brown fox jumps". Naturally the rest of the sentence can go in any way possible; therefore, the model is not sure about one specific next token (larger entropy over next tokens). However, as the sentence gets longer and we have more context like "The brown fox jumps up and down on the ground," the next tokens become more clear (smaller entropy). Also, the model is experiencing repetition which can be due to using greedy algorithm and having no tolerance for randomness. So the model goes with the exact sentence previously seen. (Also GPT2 we use here has only 12 layers <https://huggingface.co/gpt2/blob/main/config.json> So we cannot expect much from it. Larger models would generate more meaningful next tokens.)

▼ Question 1.3 (10 points)

Run the cells below to visualize the word probabilities when using different beam sizes.

```
1 input_text = 'Once upon a time, in a barn near a farm house,'
2 num_decode_steps = 20
3 model.cuda()
4
5 beam_size_list = [1, 2, 3, 4, 5]
6 output_list = []
7 probs_list = []
8 for bm in beam_size_list:
```



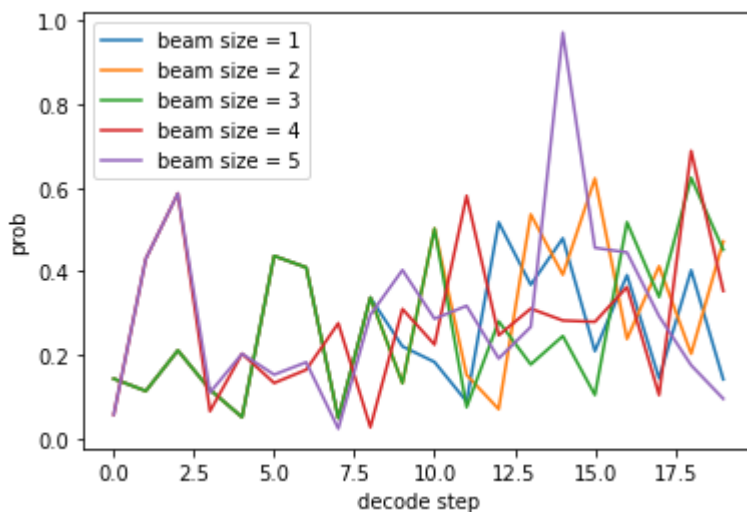
```

9     beam_output = run_beam_search(model, tokenizer, input_text, num_beams=bm, num_decode_s
10    output_list.append(beam_output)
11    probs = beam_output['token_scores'][0, 1:].exp()
12    probs_list.append((bm, probs))

1    print('Visualization with plot:')
2    fig, ax = plt.subplots()
3    for bm, probs in probs_list:
4        plt.plot(range(len(probs)), probs, label=f'beam size = {bm}')
5    plt.xlabel('decode step')
6    plt.ylabel('prob')
7    plt.legend(loc='best')
8    plt.show()
9
10   print('Model predictions:')
11   for bm, beam_output in zip(beam_size_list, output_list):
12       tokens = beam_output['output_ids'][0]
13       print(bm, beam_output['beam_scores'][0].item() / tokens.shape[-1], tokenizer.decode(tc

```

Visualization with plot:



Model predictions:

```

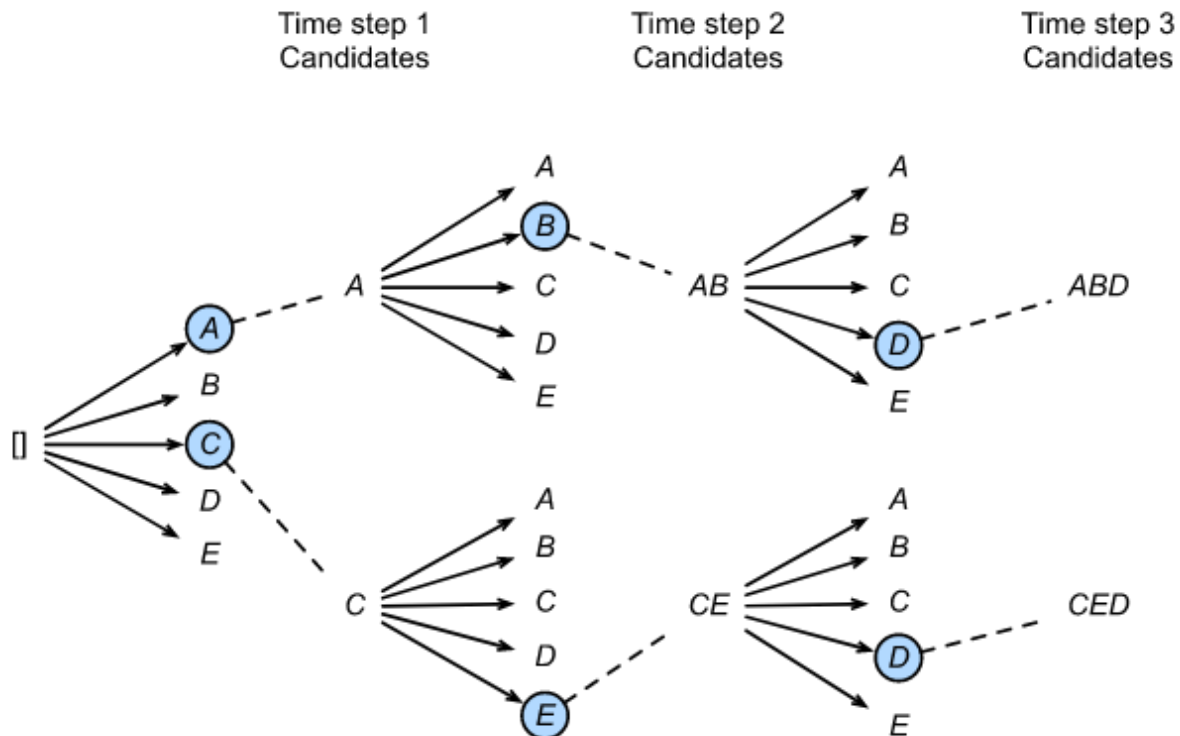
1 -0.9706196640477036 Once upon a time, in a barn near a farm house, a young boy was pl
2 -0.9286184021920869 Once upon a time, in a barn near a farm house, a young boy was pl
3 -0.9597580071651575 Once upon a time, in a barn near a farm house, a young boy was pl
4 -0.920514540238814 Once upon a time, in a barn near a farm house, there was a young g
5 -0.9058765064586293 Once upon a time, in a barn near a farm house, there was a man wh

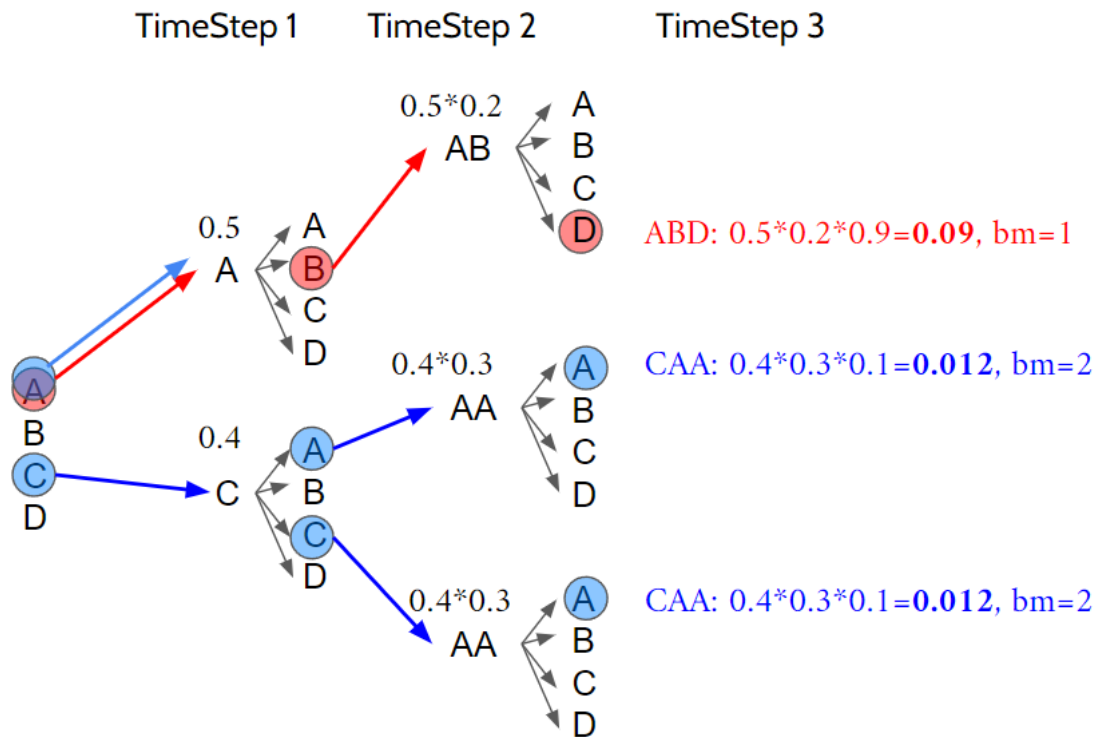
```

The Model Predictions section above includes the average cumulative log probability of each sequence. Does higher beam size always guarantee a higher probability final sequence? Why or why not?

WRITE YOUR ANSWER HERE IN A FEW SENTENCES

Beam search is an optimization of Best-First Search only going through specific beam size of all branches in each time step. As seen above, higher beam size does not always guarantee a higher probability final sequence. This is because beam search is not optimal, meaning there is no guarantee that it will find the best solution. Also, with two different beam widths, the algorithm might select specific branches and discard others in each of them. There is no guarantee that larger beam width contains the branches analyzed in a smaller one. So, it might discard a branch of $bw=1$ early; however, in the end that branch would become the highest general probability.





▼ Question 1.4 (15 points)

Beam search often results in repetition in the predicted tokens. In the following cell we pass a score processor called `WordBlock` to `run_beam_search`. At each time step, it reduces the probability for any previously seen word so that it is not generated again.

Run the cells to see how the output of beam search changes with and without using `WordBlock`.

```
1 class WordBlock:
2     def __call__(self, input_ids, scores):
3         for batch_idx in range(input_ids.shape[0]):
4             for x in input_ids[batch_idx].tolist():
5                 scores[batch_idx, x] = -1e9
6         return scores

1 input_text = 'Once upon a time, in a barn near a farm house,'
2 num_beams = 1
3
4 print('Beam Search')
5 beam_output = run_beam_search(model, tokenizer, input_text, num_beams=num_beams, num_dec
6 print(tokenizer.decode(beam_output['output_ids'][0], skip_special_tokens=True))
7
8 print('Beam Search w/ Word Block')
9 beam_output = run_beam_search(model, tokenizer, input_text, num_beams=num_beams, num_dec
```

```
10 print(tokenizer.decode(beam_output['output_ids'][0], skip_special_tokens=True))
```

```
11
```

```
Beam Search
```

```
Once upon a time, in a barn near a farm house, a young boy was playing with a stick. He
```

```
Beam Search w/ Word Block
```

```
Once upon a time, in a barn near a farm house, the young girl was playing with her fath
```

Is WordBlock a practical way to prevent repetition in beam search? What (if anything) could go wrong when using WordBlock?

WRITE YOUR ANSWER HERE IN A FEW SENTENCES

The biggest issue with wordblock is that after a while, there will be no words left to generate! This might be fixed by applying wordblock in a limited window before the next token, not all the previous tokens (Although this way the model can still repeat itself in longer sentences). Also, most frequent words such as "the", "a", and "and" will not be allowed to be used more than once which is restricting for the model to even generate new sentences. These issues can be partially fixed by the beamblock and blocking n-grams instead of unigrams but not completely.

▼ Question 1.5 (20 points)

Use the previous WordBlock example to write a new score processor called BeamBlock. Instead of uni-grams, your implementation should prevent tri-grams from appearing more than once in the sequence.

Note: This technique is called "beam blocking" and is described [here](#) (section 2.5). Also, for this assignment you do not need to re-normalize your output distribution after masking values, although typically re-normalization is done.

Write your code in the indicated section in the below cell.

```
1 class BeamBlock:
2     def __call__(self, input_ids, scores):
3         for batch_idx in range(input_ids.shape[0]):
4             # WRITE YOUR CODE HERE!
5             trigram = []
6             for x in input_ids[batch_idx].tolist():
7                 trigram.append(x)
8                 if len(trigram) > 3:
9                     trigram.pop(0)
10                if len(trigram) == 3 and torch.equal(torch.tensor(trigram[:-1]), input_i
11                print("Preventing: ", tokenizer.decode(input_ids[batch_idx, -2:]), '
```

```

12             scores[batch_idx, x] = -1e9
13         return scores

1 input_text = 'Once upon a time, in a barn near a farm house,'
2 num_beams = 1
3
4 print('Beam Search')
5 beam_output = run_beam_search(model, tokenizer, input_text, num_beams=num_beams, num_dec
6 print(tokenizer.decode(beam_output['output_ids'][0], skip_special_tokens=True))
7
8 print('Beam Search w/ Beam Block')
9 beam_output = run_beam_search(model, tokenizer, input_text, num_beams=num_beams, num_dec
10 print(tokenizer.decode(beam_output['output_ids'][0], skip_special_tokens=True))
11

```

Beam Search

Once upon a time, in a barn near a farm house, a young boy was playing with a stick. He

Beam Search w/ Beam Block

Preventing: was playing -> with

Preventing: boy was -> playing

Preventing: boy was -> playing

Preventing: boy was -> trying

Preventing: the stick -> ,

Preventing: , and -> the

Once upon a time, in a barn near a farm house, a young boy was playing with a stick. He

▼ Part 2. Language Model Fine-tuning

Now, we'll switch over to *fine-tuning* a pretrained language model. For this task, we'll use data from the [Conversational Question Answering dataset \(CoQA\)](#). The CoQA dataset includes tuples of (story text, question, answers), and we'll only be using the story text which come from various sources including children's stories, news passages, and wikipedia.

Run the below cell to set some stuff up.

```

1 import logging
2 import math
3 import os
4 import sys
5 from dataclasses import dataclass, field
6 from itertools import chain
7 from typing import Optional
8
9 import datasets
10 from datasets import load_dataset, load_metric
11
12 import transformers

```

```

13 from transformers import (
14     CONFIG_MAPPING,
15     MODEL_FOR_CAUSAL_LM_MAPPING,
16     AutoConfig,
17     AutoModelForCausalLM,
18     AutoTokenizer,
19     HfArgumentParser,
20     Trainer,
21     TrainingArguments,
22     default_data_collator,
23     is_torch_tpu_available,
24     set_seed,
25 )
26 from transformers.testing_utils import CaptureLogger
27 from transformers.trainer_utils import get_last_checkpoint
28 from transformers.utils import check_min_version
29 from transformers.utils.versions import require_version
30
31 import copy
32 import torch
33
34 from tqdm import tqdm
35 import collections
36 import numpy as np

1  MODEL_CONFIG_CLASSES = list(MODEL_FOR_CAUSAL_LM_MAPPING.keys())
2  MODEL_TYPES = tuple(conf.model_type for conf in MODEL_CONFIG_CLASSES)
3
4  @dataclass
5  class ModelArguments:
6      """
7      Arguments pertaining to which model/config/tokenizer we are going to fine-tune, or t
8      """
9      model_name_or_path: Optional[str] = field(
10         default=None,
11         metadata={
12             "help": "The model checkpoint for weights initialization."
13             "Don't set if you want to train a model from scratch."
14         },
15     )
16     model_type: Optional[str] = field(
17         default=None,
18         metadata={"help": "If training from scratch, pass a model type from the list: "
19     )
20     config_overrides: Optional[str] = field(
21         default=None,
22         metadata={
23             "help": "Override some existing default config settings when a model is trai
24             "n_embd=10,resid_pdrop=0.2,scale_attn_weights=false,summary_type=cls_index"
25         },
26     )

```

```

27     config_name: Optional[str] = field(
28         default=None, metadata={"help": "Pretrained config name or path if not the same
29     )
30     tokenizer_name: Optional[str] = field(
31         default=None, metadata={"help": "Pretrained tokenizer name or path if not the sa
32     )
33     cache_dir: Optional[str] = field(
34         default=None,
35         metadata={"help": "Where do you want to store the pretrained models downloaded f
36     )
37     use_fast_tokenizer: bool = field(
38         default=True,
39         metadata={"help": "Whether to use one of the fast tokenizer (backed by the token
40     )
41     model_revision: str = field(
42         default="main",
43         metadata={"help": "The specific model version to use (can be a branch name, tag
44     )
45     use_auth_token: bool = field(
46         default=False,
47         metadata={
48             "help": "Will use the token generated when running `transformers-cli login`
49             "with private models).",
50         },
51     )
52
53     def __post_init__(self):
54         if self.config_overrides is not None and (self.config_name is not None or self.n
55             raise ValueError(
56                 "--config_overrides can't be used in combination with --config_name or -
57         )
58
59
60 @dataclass
61 class DataTrainingArguments:
62     """
63     Arguments pertaining to what data we are going to input our model for training and e
64     """
65     dataset_name: Optional[str] = field(
66         default=None, metadata={"help": "The name of the dataset to use (via the dataset
67     )
68     dataset_config_name: Optional[str] = field(
69         default=None, metadata={"help": "The configuration name of the dataset to use (\
70     )
71     train_file: Optional[str] = field(default=None, metadata={"help": "The input trainin
72     validation_file: Optional[str] = field(
73         default=None,
74         metadata={"help": "An optional input evaluation data file to evaluate the people
75     )
76     max_train_samples: Optional[int] = field(
77         default=None,
78         metadata={

```

```

79         "help": "For debugging purposes or quicker training, truncate the number of
80         "value if set."
81     },
82 )
83 max_eval_samples: Optional[int] = field(
84     default=None,
85     metadata={
86         "help": "For debugging purposes or quicker training, truncate the number of
87         "value if set."
88     },
89 )
90
91 block_size: Optional[int] = field(
92     default=None,
93     metadata={
94         "help": "Optional input sequence length after tokenization. "
95         "The training dataset will be truncated in block of this size for training.
96         "Default to the model max input length for single sentence inputs (take into
97     },
98 )
99 overwrite_cache: bool = field(
100     default=False, metadata={"help": "Overwrite the cached training and evaluation s
101 )
102 validation_split_percentage: Optional[int] = field(
103     default=5,
104     metadata={
105         "help": "The percentage of the train set used as validation set in case ther
106     },
107 )
108 preprocessing_num_workers: Optional[int] = field(
109     default=None,
110     metadata={"help": "The number of processes to use for the preprocessing."},
111 )
112 keep_linebreaks: bool = field(
113     default=True, metadata={"help": "Whether to keep line breaks when using TXT file
114 )
115
116 def __post_init__(self):
117     if self.dataset_name is None and self.train_file is None and self.validation_file
118         raise ValueError("Need either a dataset name or a training/validation file."
119     else:
120         if self.train_file is not None:
121             extension = self.train_file.split(".")[-1]
122             assert extension in ["csv", "json", "txt"], "`train_file` should be a cs
123         if self.validation_file is not None:
124             extension = self.validation_file.split(".")[-1]
125             assert extension in ["csv", "json", "txt"], "`validation_file` should be
126

```

1 # Copied from huggingface examples.

2 #


```

1  #
2  # Modified to include the following features:
3  # - Run as a command using arguments pass as a dictionary.
4  # - Returns the model before and after fine-tuning.
5
6
7  #!/usr/bin/env python
8  # coding=utf-8
9  # Copyright 2020 The HuggingFace Inc. team. All rights reserved.
10 #
11 # Licensed under the Apache License, Version 2.0 (the "License");
12 # you may not use this file except in compliance with the License.
13 # You may obtain a copy of the License at
14 #
15 # http://www.apache.org/licenses/LICENSE-2.0
16 #
17 # Unless required by applicable law or agreed to in writing, software
18 # distributed under the License is distributed on an "AS IS" BASIS,
19 # WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
20 # See the License for the specific language governing permissions and
21 # limitations under the License.
22 """
23 Fine-tuning the library models for causal language modeling (GPT, GPT-2, CTRL, ...) on a
24 Here is the full list of checkpoints on the hub that can be fine-tuned by this script:
25 https://huggingface.co/models?filter=text-generation
26 """
27 # You can also adapt this script on your own causal language modeling task. Pointers for
28 # require_version("datasets>=1.8.0", "To fix: pip install -r examples/pytorch/language-n
29
30 logger = logging.getLogger(__name__)
31
32 MODEL_CONFIG_CLASSES = list(MODEL_FOR_CAUSAL_LM_MAPPING.keys())
33 MODEL_TYPES = tuple(conf.model_type for conf in MODEL_CONFIG_CLASSES)
34
35
36 def run_clm(args_as_dict, debug_state={}):
37
38     # See all possible arguments in src/transformers/training_args.py
39     # or by passing the --help flag to this script.
40     # We now keep distinct sets of args, for a cleaner separation of concerns.
41
42     parser = HfArgumentParser((ModelArguments, DataTrainingArguments, TrainingArguments))
43     model_args, data_args, training_args = parser.parse_dict(args_as_dict)
44
45     logging.basicConfig(
46         format="%(asctime)s - %(levelname)s - %(name)s - %(message)s",
47         datefmt="%m/%d/%Y %H:%M:%S",
48         handlers=[logging.StreamHandler(sys.stdout)],
49     )
50
51     log_level = training_args.get_process_log_level()
52     logger.setLevel(log_level)
53     datasets.utils.logging.set_verbosity(log_level)

```

```

54 transformers.utils.logging.set_verbosity(log_level)
55 transformers.utils.logging.enable_default_handler()
56 transformers.utils.logging.enable_explicit_format()
57
58 logger.warning(
59     f"Process rank: {training_args.local_rank}, device: {training_args.device}, n_g
60     + f"distributed training: {bool(training_args.local_rank != -1)}, 16-bits traini
61 )
62 logger.info(f"Training/evaluation parameters {training_args}")
63
64 # Detecting last checkpoint.
65 last_checkpoint = None
66 if os.path.isdir(training_args.output_dir) and training_args.do_train and not traini
67     last_checkpoint = get_last_checkpoint(training_args.output_dir)
68     if last_checkpoint is None and len(os.listdir(training_args.output_dir)) > 0:
69         raise ValueError(
70             f"Output directory ({training_args.output_dir}) already exists and is no
71             "Use --overwrite_output_dir to overcome."
72         )
73     elif last_checkpoint is not None and training_args.resume_from_checkpoint is Nor
74         logger.info(
75             f"Checkpoint detected, resuming training at {last_checkpoint}. To avoid
76             "the `--output_dir` or add `--overwrite_output_dir` to train from scratc
77         )
78
79 # Set seed before initializing model.
80 set_seed(training_args.seed)
81
82 data_files = {}
83 dataset_args = {}
84 if data_args.train_file is not None:
85     data_files["train"] = data_args.train_file
86 if data_args.validation_file is not None:
87     data_files["validation"] = data_args.validation_file
88 raw_datasets = load_dataset('json', data_files=data_files, cache_dir=model_args.cach
89
90 # See more about loading any type of standard or custom dataset (from files, python
91 # https://huggingface.co/docs/datasets/loading\_datasets.html).
92
93 # Load pretrained model and tokenizer
94 #
95 # Distributed training:
96 # The .from_pretrained methods guarantee that only one local process can concurrentl
97 # download model & vocab.
98
99 config_kwargs = {
100     "cache_dir": model_args.cache_dir,
101     "revision": model_args.model_revision,
102     "use_auth_token": True if model_args.use_auth_token else None,
103 }
104 if model_args.config_name:
105     config = AutoConfig.from_pretrained(model_args.config_name, **config_kwargs)

```

```

106 elif model_args.model_name_or_path:
107     config = AutoConfig.from_pretrained(model_args.model_name_or_path, **config_kwar
108 else:
109     config = CONFIG_MAPPING[model_args.model_type]()
110     logger.warning("You are instantiating a new config instance from scratch.")
111     if model_args.config_overrides is not None:
112         logger.info(f"Overriding config: {model_args.config_overrides}")
113         config.update_from_string(model_args.config_overrides)
114         logger.info(f"New config: {config}")
115
116 tokenizer_kwargs = {
117     "cache_dir": model_args.cache_dir,
118     "use_fast": model_args.use_fast_tokenizer,
119     "revision": model_args.model_revision,
120     "use_auth_token": True if model_args.use_auth_token else None,
121 }
122 if model_args.tokenizer_name:
123     tokenizer = AutoTokenizer.from_pretrained(model_args.tokenizer_name, **tokenizer
124 elif model_args.model_name_or_path:
125     tokenizer = AutoTokenizer.from_pretrained(model_args.model_name_or_path, **token
126 else:
127     raise ValueError(
128         "You are instantiating a new tokenizer from scratch. This is not supported b
129         "You can do it from another script, save it, and load it from here, using --
130     )
131
132 debug_state['tokenizer'] = tokenizer
133
134 if model_args.model_name_or_path:
135     model = AutoModelForCausalLM.from_pretrained(
136         model_args.model_name_or_path,
137         from_tf=bool(".ckpt" in model_args.model_name_or_path),
138         config=config,
139         cache_dir=model_args.cache_dir,
140         revision=model_args.model_revision,
141         use_auth_token=True if model_args.use_auth_token else None,
142     )
143 else:
144     model = AutoModelForCausalLM.from_config(config)
145     n_params = sum(dict((p.data_ptr(), p.numel()) for p in model.parameters()).value
146     logger.info(f"Training new model from scratch - Total size={n_params/2**20:.2f}M
147
148 model.resize_token_embeddings(len(tokenizer))
149
150 model_before_finetuning = debug_state["model_before_finetuning"] = copy.deepcopy(moc
151
152 # Preprocessing the datasets.
153 # First we tokenize all the texts.
154 if training_args.do_train:
155     column_names = raw_datasets["train"].column_names
156 else:

```

```

157     column_names = raw_datasets["validation"].column_names
158     text_column_name = "story" if "story" in column_names else column_names[0]
159
160     if args_as_dict.get('text_column_name', None) is not None:
161         text_column_name = args_as_dict['text_column_name']
162
163     # since this will be pickled to avoid _LazyModule error in Hasher force logger loadi
164     tok_logger = transformers.utils.logging.get_logger("transformers.tokenization_utils_
165
166     def tokenize_function(examples):
167         with CaptureLogger(tok_logger) as cl:
168             output = tokenizer(examples[text_column_name])
169             # clm input could be much much longer than block_size
170             if "Token indices sequence length is longer than the" in cl.out:
171                 tok_logger.warning(
172                     "^^^^^^^^^^^^^^^^^^^^ Please ignore the warning above - this long input will
173                 )
174         return output
175
176     with training_args.main_process_first(desc="dataset map tokenization"):
177         tokenized_datasets = raw_datasets.map(
178             tokenize_function,
179             batched=True,
180             num_proc=data_args.preprocessing_num_workers,
181             remove_columns=column_names,
182             load_from_cache_file=not data_args.overwrite_cache,
183             desc="Running tokenizer on dataset",
184         )
185
186     if data_args.block_size is None:
187         block_size = tokenizer.model_max_length
188         if block_size > 1024:
189             logger.warning(
190                 f"The tokenizer picked seems to have a very large `model_max_length` ({t
191                 "Picking 1024 instead. You can change that default value by passing --b
192             )
193             block_size = 1024
194     else:
195         if data_args.block_size > tokenizer.model_max_length:
196             logger.warning(
197                 f"The block_size passed ({data_args.block_size}) is larger than the maxi
198                 f"({tokenizer.model_max_length}). Using block_size={tokenizer.model_max_
199             )
200             block_size = min(data_args.block_size, tokenizer.model_max_length)
201
202     debug_state['block_size'] = block_size
203
204     # Main data processing function that will concatenate all texts from our dataset and
205     def group_texts(examples):
206         # Concatenate all texts.
207         concatenated_examples = {k: list(chain(*examples[k])) for k in examples.keys()}
208         total_length = len(concatenated_examples[list(examples.keys())[0]])

```

```

209     # We drop the small remainder, we could add padding if the model supported it ir
210     # customize this part to your needs.
211     if total_length >= block_size:
212         total_length = (total_length // block_size) * block_size
213     # Split by chunks of max_len.
214     result = {
215         k: [t[i : i + block_size] for i in range(0, total_length, block_size)]
216         for k, t in concatenated_examples.items()
217     }
218     result["labels"] = result["input_ids"].copy()
219     return result
220
221     # Note that with `batched=True`, this map processes 1,000 texts together, so group_t
222     # for each of those groups of 1,000 texts. You can adjust that batch_size here but a
223     # to preprocess.
224     #
225     # To speed up this part, we use multiprocessing. See the documentation of the map me
226     # https://huggingface.co/docs/datasets/package\_reference/main\_classes.html#datasets.
227
228     with training_args.main_process_first(desc="grouping texts together"):
229         lm_datasets = tokenized_datasets.map(
230             group_texts,
231             batched=True,
232             num_proc=data_args.preprocessing_num_workers,
233             load_from_cache_file=not data_args.overwrite_cache,
234             desc=f"Grouping texts in chunks of {block_size}",
235         )
236
237     if training_args.do_train:
238         if "train" not in tokenized_datasets:
239             raise ValueError("--do_train requires a train dataset")
240         train_dataset = lm_datasets["train"]
241         if data_args.max_train_samples is not None:
242             max_train_samples = min(len(train_dataset), data_args.max_train_samples)
243             train_dataset = train_dataset.select(range(max_train_samples))
244
245     if training_args.do_eval:
246         if "validation" not in tokenized_datasets:
247             raise ValueError("--do_eval requires a validation dataset")
248         eval_dataset = lm_datasets["validation"]
249         if data_args.max_eval_samples is not None:
250             max_eval_samples = min(len(eval_dataset), data_args.max_eval_samples)
251             eval_dataset = eval_dataset.select(range(max_eval_samples))
252
253     def preprocess_logits_for_metrics(logits, labels):
254         if isinstance(logits, tuple):
255             # Depending on the model and config, logits may contain extra tensors,
256             # like past_key_values, but logits always come first
257             logits = logits[0]
258         return logits.argmax(dim=-1)
259
260     metric = load_metric("accuracy")

```

```

260 metric = load_metric(accuracy)
261
262 def compute_metrics(eval_preds):
263     preds, labels = eval_preds
264     # preds have the same shape as the labels, after the argmax(-1) has been cal
265     # by preprocess_logits_for_metrics but we need to shift the labels
266     labels = labels[:, 1:].reshape(-1)
267     preds = preds[:, :-1].reshape(-1)
268     return metric.compute(predictions=preds, references=labels)
269
270 # Initialize our Trainer
271 trainer = Trainer(
272     model=model,
273     args=training_args,
274     train_dataset=train_dataset if training_args.do_train else None,
275     eval_dataset=eval_dataset if training_args.do_eval else None,
276     tokenizer=tokenizer,
277     # Data collator will default to DataCollatorWithPadding, so we change it.
278     data_collator=default_data_collator,
279     compute_metrics=compute_metrics if training_args.do_eval and not is_torch_tpu_av
280     preprocess_logits_for_metrics=preprocess_logits_for_metrics
281     if training_args.do_eval and not is_torch_tpu_available()
282     else None,
283 )
284
285 # Training
286 model_after_finetuning = debug_state["model_after_finetuning"] = None
287 if training_args.do_train:
288     checkpoint = None
289     if training_args.resume_from_checkpoint is not None:
290         checkpoint = training_args.resume_from_checkpoint
291     elif last_checkpoint is not None:
292         checkpoint = last_checkpoint
293     train_result = trainer.train(resume_from_checkpoint=checkpoint)
294     trainer.save_model() # Saves the tokenizer too for easy upload
295
296     metrics = train_result.metrics
297
298     max_train_samples = (
299         data_args.max_train_samples if data_args.max_train_samples is not None else
300     )
301     metrics["train_samples"] = min(max_train_samples, len(train_dataset))
302
303     trainer.log_metrics("train", metrics)
304     trainer.save_metrics("train", metrics)
305     trainer.save_state()
306
307     model_after_finetuning = debug_state["model_after_finetuning"] = model
308
309 # Evaluation
310 if training_args.do_eval:
311     logger.info("*** Evaluate ***")

```

```

312
313     metrics = trainer.evaluate()
314
315     max_eval_samples = data_args.max_eval_samples if data_args.max_eval_samples is not None else
316     metrics["eval_samples"] = min(max_eval_samples, len(eval_dataset))
317     try:
318         perplexity = math.exp(metrics["eval_loss"])
319     except OverflowError:
320         perplexity = float("inf")
321     metrics["perplexity"] = perplexity
322
323     trainer.log_metrics("eval", metrics)
324     trainer.save_metrics("eval", metrics)
325
326     kwargs = {"finetuned_from": model_args.model_name_or_path, "tasks": "text-generation"}
327     if data_args.dataset_name is not None:
328         kwargs["dataset_tags"] = data_args.dataset_name
329         if data_args.dataset_config_name is not None:
330             kwargs["dataset_args"] = data_args.dataset_config_name
331             kwargs["dataset"] = f"{data_args.dataset_name} {data_args.dataset_config_name}"
332         else:
333             kwargs["dataset"] = data_args.dataset_name
334
335     # Should call this after `run_clm` to free up some GPU memory.
336     # Some GPU memory will still be reserved, so if you need to re-run
337     # fine-tuning, then you may need to click "Runtime -> Restart Runtime", although
338     # this will reset all previously run cells.
339     model_before_finetuning.cpu()
340     model_after_finetuning.cpu()
341     torch.cuda.empty_cache()
342
343     return model_before_finetuning, model_after_finetuning

```

```

1 def compute_rouge(model, tokenizer, dataset, n=3):
2
3     def count_ngrams(tokens, n):
4         c = collections.Counter()
5         for size in range(1, n + 1):
6             for end in range(size, len(tokens) + 1):
7                 ngram = tuple(tokens[end - size:end])
8                 c[ngram] += 1
9         return c
10
11     def rouge(gold, pred, n):
12         gold_c = count_ngrams(gold, n)
13         pred_c = count_ngrams(pred, n)
14         overlap = sum([pred_c[ngram] for ngram in gold_c.keys()])
15         total = sum(gold_c.values())
16         return overlap / total
17

```

```

18 with torch.inference_mode():
19     m = []
20     for p1, p2 in tqdm(dataset, desc=f'Compute ROGUE-{n}'):
21         # TODO: Does this include the correct values for beam search?
22         beam_output = run_beam_search(
23             model,
24             tokenizer,
25             p1,
26             num_beams=3,
27             num_decode_steps=32)
28         pred = tokenizer.decode(beam_output['output_ids'][0], skip_special_tokens=True).sp
29         pred_ids = tokenizer(pred, return_tensors="pt")['input_ids'][0].tolist()
30         # p1_tensor = tokenizer(p1, return_tensors="pt")['input_ids']
31         gold_ids = tokenizer(p2, return_tensors="pt")['input_ids'][0].tolist()
32         m.append(rouge(gold_ids, pred_ids, n))
33
34     return np.mean(m)
35
36
37 def compute_perplexity(model, tokenizer, dataset):
38
39     with torch.inference_mode():
40         n = 0
41         m = []
42         for p1, p2 in tqdm(dataset, desc='Compute Perplexity'):
43             p1_tensor = tokenizer(p1, return_tensors="pt")['input_ids']
44             p2_tensor = tokenizer(p2, return_tensors="pt")['input_ids']
45             input_ids = torch.cat([p1_tensor, p2_tensor], 1).to(model.device)
46             target = input_ids.clone()
47             target[:, :p1_tensor.shape[1]] = -100
48             target_length = p2_tensor.shape[1]
49             n += target_length
50
51             nll = model(input_ids=input_ids, labels=target)[0] * target_length
52             m.append(nll)
53
54     return torch.exp(torch.cat([x.view(1) for x in m], 0).sum() / n)
55
56
57 def preprocess_coqa(dataset):
58     new_dataset = []
59     skipped = 0
60     for text in dataset:
61         parts = text.split('. ', 2)
62         if len(parts) <= 1:
63             skipped += 1
64             continue
65         p1 = parts[0].strip() + '.'
66         p2 = parts[1].strip() + '.'
67         new_dataset.append((p1, p2))

```


17

▼ Question 2.1 (15 points)

Run the cell below, which does the following steps:

- Fine-tune GPT-2 on the story text from the CoQA dataset.
- Preprocess the CoQA dataset into "sentence pairs". These pairs are created by finding the first and second sentence from each story passage in the validation data.
- Evaluate the language models from before and after fine-tuning using two different metrics: perplexity and ROUGE-3.

Important Notes:

- For training, the full story passages are used and can be many sentences long. For evaluation, only the sentence pairs are used. Both evaluation metrics are only evaluated on the second sentence.
- For perplexity, we use teacher forcing. For ROUGE-3 we use the first sentence as a prefix and generate a second sentence using beam search (beam size of 3 and generating a fixed amount of 32 tokens).

```
1 # Fine-tune GPT-2
2
3 config = {
4     'model_name_or_path': 'gpt2',
5     'train_file': 'coqa-train.json',
6     'validation_file': 'coqa-dev.json',
7     'text_column_name': 'story',
8     'per_device_train_batch_size': 8,
9     'per_device_eval_batch_size': 8,
10    'gradient_accumulation_steps': 1,
11    'learning_rate': 5e-5,
12    'block_size': 256,
13    'max_train_samples': 1024,
14    'num_train_epochs': 1,
15    'do_train': True,
16    'do_eval': False,
17    'output_dir': './tmp',
18    'overwrite_output_dir': True,
19    'log_level': 'warning' # Set to `info` or `debug` for additional logging.
20 }
21
22 # If preferred, can use these arguments in the config instead of `train_file`
23 # and `validation_file`, but sometimes Google Colab's IP gets throttled.
24 # 'dataset_name': 'coqa',
25 # 'dataset_config_name': 'default',
26
```

```
27 model_before_finetuning, model_after_finetuning = run_clm(config)
28 print('LM finetuning finished!')
29
30 # Preprocess the CoQA dataset into sentence pairs for evaluation.
31
32 dataset = load_dataset('json', data_files={'validation':'coqa-dev.json'}, field='data')[
33 new_dataset = preprocess_coqa(dataset)
34
35 print('Preprocessing finished!')
36 print(f'...found {len(new_dataset)} instances.')
37 print(f'...sample instance: {new_dataset[0]}')
38
39 # Run evaluation.
40
41 model_before_finetuning.cuda()
42 model_before_finetuning.eval()
43 model_after_finetuning.cuda()
44 model_after_finetuning.eval()
45
46 print('Running evaluation...')
47
48 before_ppl = compute_perplexity(model_before_finetuning, tokenizer, new_dataset).item()
49 after_ppl = compute_perplexity(model_after_finetuning, tokenizer, new_dataset).item()
50 print(f'\n\nPerplexity before_finetune = {before_ppl:.3f}, after_finetune = {after_ppl:.3f}')
51
52 before_rouge = compute_rouge(model_before_finetuning, tokenizer, new_dataset)
53 after_rouge = compute_rouge(model_after_finetuning, tokenizer, new_dataset)
54 print(f'\n\nROUGE-3 before_finetune = {before_rouge:.3f}, after_finetune = {after_rouge:.3f}')
55
56 print('Evaluation finished!')
57
```

```

WARNING:__main__:Process rank: -1, device: cuda:0, n_gpu: 1distributed training: False,
WARNING:datasets.builder:Using custom data configuration default-4a9a50cb04c53708
Downloading and preparing dataset json/default to /root/.cache/huggingface/datasets/jso
Downloading data files: 100% 2/2 [00:00<00:00, 57.29it/s]

Extracting data files: 100% 2/2 [00:00<00:00, 75.09it/s]
Dataset json downloaded and prepared to /root/.cache/huggingface/datasets/json/default-
100% 2/2 [00:00<00:00, 45.96it/s]

Downloading: 100% 1.29M/1.29M [00:01<00:00, 928kB/s]
WARNING:datasets.fingerprint:Parameter 'function'=<function run_clm.<locals>.tokenize_f
Running tokenizer on dataset: 8/8 [00:06<00:00,
100% 1.49ba/s]
[WARNING|tokenization_utils_base.py:3397] 2022-12-03 15:44:26,772 >> Token indices sequ
[WARNING|<ipython-input-43-6b0a8fc7b51e>:171] 2022-12-03 15:44:26,774 >> ^^^^^^^^^^^^^^^^^
Running tokenizer on dataset: 1/1 [00:00<00:00,
100% 2.11ba/s]

Grouping texts in chunks of 256: 8/8 [00:02<00:00,
100% 2.64ba/s]

Grouping texts in chunks of 256: 1/1 [00:00<00:00,
100% 4.42ba/s]

/usr/local/lib/python3.8/dist-packages/transformers/optimization.py:306: FutureWarning:
warnings.warn(
[128/128 01:14, Epoch 1/1]

```

Step Training Loss

```
***** train metrics *****
```

```

epoch                =      1.0
total_flos           = 124593GF
train_loss            =    3.4296
train_runtime        = 0:01:15.59
train_samples         =    1024
train_samples_per_second =    13.546
train_steps_per_second  =    1.693

```

Has language model fine-tuning improved GPT-2 performance on the story text for the CoQA dataset? Is perplexity or ROUGE a better metric for measuring this?

```

Downloading data files: 100% 1/1 [00:00<00:00, 29.11it/s]

```

WRITE YOUR ANSWER HERE IN A FEW SENTENCES

Perplexity before_finetune = 34.129, after_finetune = 30.469

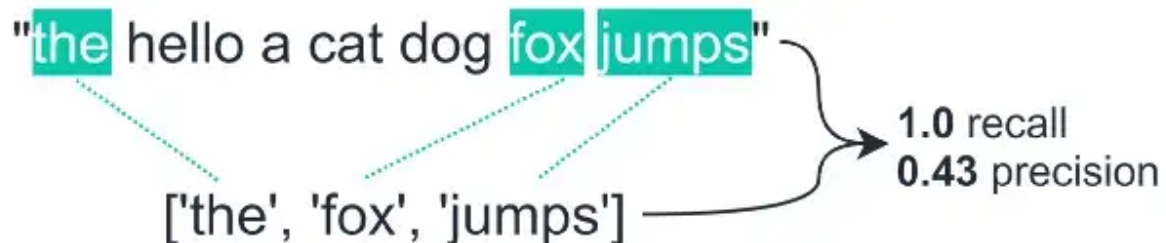
ROUGE-3 before_finetune = 0.086, after_finetune = 0.083

The language model fine-tuning improved GPT-2 performance on the story text for the CoQA dataset. The perplexity decreased from 34 to 30 (average branching factor

reduced) and ROUGE-3 did not change much. (0.086->0.083)

Perplexity is a better metric for measuring this. ROUGE is more used for evaluating automatic summarization and machine translation, where the generated text should include specific words and phrases to best express the prior input text. On the other hand, in our case of generating second sentence of a story given the previous sentence, this does not hold. The second sentence can have many possibilities and all can be correct even without much sharing ngrams with the real sentence.

Therefore, in this case and also evaluating language models in general for their generation capabilities, perplexity best represents their power. It best captures whether the model has properly fit the data or not, which in this case is the context of stories.



$$2 * \frac{0.43 * 1.0}{0.43 + 1.0} = 0.6$$

60% f1 score

Perplexity

- Does the model fit the data?
 - A good model will give a high probability to a real sentence
- Perplexity
 - Average branching factor in predicting the next word
 - Lower is better (lower perplexity -> higher probability)
 - N = number of words

$$Per = \sqrt[N]{\frac{1}{P(w_1 w_2 \dots w_N)}}$$

We're done with GPT-2 for now, so move the trained models back to CPU.

```
1 # Should call this after `run_clm` to free up some GPU memory.
2 # Some GPU memory will still be reserved, so if you need to re-run
3 # fine-tuning, then you may need to click "Runtime -> Restart Runtime", although
4 # this will reset all previously run cells.
5 model_before_finetuning.cpu()
6 model_after_finetuning.cpu()
7 torch.cuda.empty_cache()
```

▼ Part 3: Data Augmentation via Backtranslation

The last part of this homework involves data augmentation of an NLP classifier via backtranslation. Now run the below cell to set up some fine-tuning code.

```
1 import logging
2 import os
3 import random
4 import sys
5 from dataclasses import dataclass, field
6 from typing import Optional
7
8 import datasets
```

```
9 import numpy as np
10 from datasets import load_dataset, load_metric
11
12 import transformers
13 from transformers import (
14     AutoConfig,
15     AutoModelForSequenceClassification,
16     AutoTokenizer,
17     DataCollatorWithPadding,
18     EvalPrediction,
19     HfArgumentParser,
20     PretrainedConfig,
21     Trainer,
22     TrainingArguments,
23     default_data_collator,
24     set_seed,
25 )
26 from transformers.trainer_utils import get_last_checkpoint
27 from transformers.utils import check_min_version
28 from transformers.utils.versions import require_version
29
30 from transformers import glue_processors
31 from transformers.data.processors.utils import InputExample
32 from langdetect import detect

1 #!/usr/bin/env python
2 # coding=utf-8
3 # Copyright 2020 The HuggingFace Inc. team. All rights reserved.
4 #
5 # Licensed under the Apache License, Version 2.0 (the "License");
6 # you may not use this file except in compliance with the License.
7 # You may obtain a copy of the License at
8 #
9 #     http://www.apache.org/licenses/LICENSE-2.0
10 #
11 # Unless required by applicable law or agreed to in writing, software
12 # distributed under the License is distributed on an "AS IS" BASIS,
13 # WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
14 # See the License for the specific language governing permissions and
15 # limitations under the License.
16 """ Finetuning the library models for sequence classification on GLUE."""
17 # You can also adapt this script on your own text classification task. Pointers for this
18
19 # Will error if the minimal version of Transformers is not installed. Remove at your own
20 # check_min_version("4.18.0.dev0")
21
22 # require_version("datasets>=1.8.0", "To fix: pip install -r examples/pytorch/text-class
23
24 task_to_keys = {
25     "cola": ("sentence", None),
26     "mnli": ("premise", "hypothesis"),
```

```

27     "mrpc": ("sentence1", "sentence2"),
28     "qnli": ("question", "sentence"),
29     "qqp": ("question1", "question2"),
30     "rte": ("sentence1", "sentence2"),
31     "sst2": ("sentence", None),
32     "stsb": ("sentence1", "sentence2"),
33     "wnli": ("sentence1", "sentence2"),
34 }
35
36 logger = logging.getLogger(__name__)
37
38
39 @dataclass
40 class DataTrainingArguments:
41     """
42     Arguments pertaining to what data we are going to input our model for training and e
43
44     Using `HfArgumentParser` we can turn this class
45     into argparse arguments to be able to specify them on
46     the command line.
47     """
48
49     task_name: Optional[str] = field(
50         default=None,
51         metadata={"help": "The name of the task to train on: " + ", ".join(task_to_keys.
52     )
53     dataset_name: Optional[str] = field(
54         default=None, metadata={"help": "The name of the dataset to use (via the dataset
55     )
56     dataset_config_name: Optional[str] = field(
57         default=None, metadata={"help": "The configuration name of the dataset to use (v
58     )
59     max_seq_length: int = field(
60         default=128,
61         metadata={
62             "help": "The maximum total input sequence length after tokenization. Sequenc
63             "than this will be truncated, sequences shorter will be padded."
64         },
65     )
66     overwrite_cache: bool = field(
67         default=False, metadata={"help": "Overwrite the cached preprocessed datasets or
68     )
69     pad_to_max_length: bool = field(
70         default=True,
71         metadata={
72             "help": "Whether to pad all samples to `max_seq_length`. "
73             "If False, will pad the samples dynamically when batching to the maximum len
74         },
75     )
76     max_train_samples: Optional[int] = field(
77         default=None,

```

```

78     metadata={
79         "help": "For debugging purposes or quicker training, truncate the number of
80         "value if set."
81     },
82 )
83 max_eval_samples: Optional[int] = field(
84     default=None,
85     metadata={
86         "help": "For debugging purposes or quicker training, truncate the number of
87         "value if set."
88     },
89 )
90 max_predict_samples: Optional[int] = field(
91     default=None,
92     metadata={
93         "help": "For debugging purposes or quicker training, truncate the number of
94         "value if set."
95     },
96 )
97 train_file: Optional[str] = field(
98     default=None, metadata={"help": "A csv or a json file containing the training da
99 )
100 validation_file: Optional[str] = field(
101     default=None, metadata={"help": "A csv or a json file containing the validation
102 )
103 test_file: Optional[str] = field(default=None, metadata={"help": "A csv or a json fi
104
105 def __post_init__(self):
106     if self.task_name is not None:
107         self.task_name = self.task_name.lower()
108         if self.task_name not in task_to_keys.keys():
109             raise ValueError("Unknown task, you should pick one in " + ",".join(task
110         elif self.dataset_name is not None:
111             pass
112         elif self.train_file is None or self.validation_file is None:
113             raise ValueError("Need either a GLUE task, a training/validation file or a d
114         else:
115             train_extension = self.train_file.split(".")[1]
116             assert train_extension in ["csv", "json"], "`train_file` should be a csv or
117             validation_extension = self.validation_file.split(".")[1]
118             assert (
119                 validation_extension == train_extension
120             ), "`validation_file` should have the same extension (csv or json) as `train
121
122
123 @dataclass
124 class ModelArguments:
125     """
126     Arguments pertaining to which model/config/tokenizer we are going to fine-tune from.
127     """
128

```



```

129     model_name_or_path: str = field(
130         metadata={"help": "Path to pretrained model or model identifier from huggingface
131     )
132     config_name: Optional[str] = field(
133         default=None, metadata={"help": "Pretrained config name or path if not the same
134     )
135     tokenizer_name: Optional[str] = field(
136         default=None, metadata={"help": "Pretrained tokenizer name or path if not the sa
137     )
138     cache_dir: Optional[str] = field(
139         default=None,
140         metadata={"help": "Where do you want to store the pretrained models downloaded f
141     )
142     use_fast_tokenizer: bool = field(
143         default=True,
144         metadata={"help": "Whether to use one of the fast tokenizer (backed by the token
145     )
146     model_revision: str = field(
147         default="main",
148         metadata={"help": "The specific model version to use (can be a branch name, tag
149     )
150     use_auth_token: bool = field(
151         default=False,
152         metadata={
153             "help": "Will use the token generated when running `transformers-cli login`
154             "with private models)."
155     }.

```

```

1 def do_target_task_finetuning(args_as_dict):
2
3     # See all possible arguments in src/transformers/training_args.py
4     # or by passing the --help flag to this script.
5     # We now keep distinct sets of args, for a cleaner separation of concerns.
6
7     parser = HfArgumentParser((ModelArguments, DataTrainingArguments, TrainingArguments)
8     model_args, data_args, training_args = parser.parse_dict(args_as_dict)
9
10    # Setup logging
11    logging.basicConfig(
12        format="%(asctime)s - %(levelname)s - %(name)s - %(message)s",
13        datefmt="%m/%d/%Y %H:%M:%S",
14        handlers=[logging.StreamHandler(sys.stdout)],
15    )
16
17    log_level = training_args.get_process_log_level()
18    logger.setLevel(log_level)
19    datasets.utils.logging.set_verbosity(log_level)
20    transformers.utils.logging.set_verbosity(log_level)
21    transformers.utils.logging.enable_default_handler()
22    transformers.utils.logging.enable_explicit_format()
23

```

```

24 # Log on each process the small summary:
25 logger.warning(
26     f"Process rank: {training_args.local_rank}, device: {training_args.device}, n_gp
27     + f"distributed training: {bool(training_args.local_rank != -1)}, 16-bits traini
28 )
29 logger.info(f"Training/evaluation parameters {training_args}")
30
31 # Detecting last checkpoint.
32 last_checkpoint = None
33 if os.path.isdir(training_args.output_dir) and training_args.do_train and not traini
34     last_checkpoint = get_last_checkpoint(training_args.output_dir)
35     if last_checkpoint is None and len(os.listdir(training_args.output_dir)) > 0:
36         raise ValueError(
37             f"Output directory ({training_args.output_dir}) already exists and is no
38             "Use --overwrite_output_dir to overcome."
39         )
40     elif last_checkpoint is not None and training_args.resume_from_checkpoint is Non
41         logger.info(
42             f"Checkpoint detected, resuming training at {last_checkpoint}. To avoid
43             "the `--output_dir` or add `--overwrite_output_dir` to train from scratc
44         )
45
46 # Set seed before initializing model.
47 set_seed(training_args.seed)
48
49 # In distributed training, the load_dataset function guarantee that only one local p
50 # download the dataset.
51 if data_args.task_name is not None:
52     # Downloading and loading a dataset from the hub.
53     raw_datasets = load_dataset("glue", data_args.task_name, cache_dir=model_args.ca
54 elif data_args.dataset_name is not None:
55     # Downloading and loading a dataset from the hub.
56     raw_datasets = load_dataset(
57         data_args.dataset_name, data_args.dataset_config_name, cache_dir=model_args.
58     )
59 else:
60     # Loading a dataset from your local files.
61     # CSV/JSON training and evaluation files are needed.
62     data_files = {"train": data_args.train_file, "validation": data_args.validation_
63
64     # Get the test dataset: you can provide your own CSV/JSON test file (see below)
65     # when you use `do_predict` without specifying a GLUE benchmark task.
66     if training_args.do_predict:
67         if data_args.test_file is not None:
68             train_extension = data_args.train_file.split(".")[1]
69             test_extension = data_args.test_file.split(".")[1]
70             assert (
71                 test_extension == train_extension
72             ), "`test_file` should have the same extension (csv or json) as `train_f
73             data_files["test"] = data_args.test_file
74     else:

```

```

75         raise ValueError("Need either a GLUE task or a test file for `do_predict`")
76
77     for key in data_files.keys():
78         logger.info(f"load a local file for {key}: {data_files[key]}")
79
80     if data_args.train_file.endswith(".csv"):
81         # Loading a dataset from local csv files
82         raw_datasets = load_dataset("csv", data_files=data_files, cache_dir=model_ar
83     else:
84         # Loading a dataset from local json files
85         raw_datasets = load_dataset("json", data_files=data_files, cache_dir=model_a
86 # See more about loading any type of standard or custom dataset at
87 # https://huggingface.co/docs/datasets/loading_datasets.html.
88
89 # Labels
90 if data_args.task_name is not None:
91     is_regression = data_args.task_name == "stsb"
92     if not is_regression:
93         label_list = raw_datasets["train"].features["label"].names
94         num_labels = len(label_list)
95     else:
96         num_labels = 1
97 else:
98     # Trying to have good defaults here, don't hesitate to tweak to your needs.
99     is_regression = raw_datasets["train"].features["label"].dtype in ["float32", "fl
100 if is_regression:
101     num_labels = 1
102 else:
103     # A useful fast method:
104     # https://huggingface.co/docs/datasets/package_reference/main_classes.html#d
105     label_list = raw_datasets["train"].unique("label")
106     label_list.sort() # Let's sort it for determinism
107     num_labels = len(label_list)
108
109 # Load pretrained model and tokenizer
110 #
111 # In distributed training, the .from_pretrained methods guarantee that only one loca
112 # download model & vocab.
113 config = AutoConfig.from_pretrained(
114     model_args.config_name if model_args.config_name else model_args.model_name_or_p
115     num_labels=num_labels,
116     finetuning_task=data_args.task_name,
117     cache_dir=model_args.cache_dir,
118     revision=model_args.model_revision,
119     use_auth_token=True if model_args.use_auth_token else None,
120 )
121 tokenizer = AutoTokenizer.from_pretrained(
122     model_args.tokenizer_name if model_args.tokenizer_name else model_args.model_nam
123     cache_dir=model_args.cache_dir,
124     use_fast=model_args.use_fast_tokenizer,
125     revision=model_args.model_revision,

```

```

126     use_auth_token=True if model_args.use_auth_token else None,
127 )
128 model = AutoModelForSequenceClassification.from_pretrained(
129     model_args.model_name_or_path,
130     from_tf=bool(".ckpt" in model_args.model_name_or_path),
131     config=config,
132     cache_dir=model_args.cache_dir,
133     revision=model_args.model_revision,
134     use_auth_token=True if model_args.use_auth_token else None,
135 )
136
137 # Preprocessing the raw_datasets
138 if data_args.task_name is not None:
139     sentence1_key, sentence2_key = task_to_keys[data_args.task_name]
140 else:
141     # Again, we try to have some nice defaults but don't hesitate to tweak to your u
142     non_label_column_names = [name for name in raw_datasets["train"].column_names if
143     if "sentence1" in non_label_column_names and "sentence2" in non_label_column_nam
144         sentence1_key, sentence2_key = "sentence1", "sentence2"
145     else:
146         if len(non_label_column_names) >= 2:
147             sentence1_key, sentence2_key = non_label_column_names[:2]
148         else:
149             sentence1_key, sentence2_key = non_label_column_names[0], None
150
151 # Padding strategy
152 if data_args.pad_to_max_length:
153     padding = "max_length"
154 else:
155     # We will pad later, dynamically at batch creation, to the max sequence length i
156     padding = False
157
158 # Some models have set the order of the labels to use, so let's make sure we do use
159 label_to_id = None
160 if (
161     model.config.label2id != PretrainedConfig(num_labels=num_labels).label2id
162     and data_args.task_name is not None
163     and not is_regression
164 ):
165     # Some have all caps in their config, some don't.
166     label_name_to_id = {k.lower(): v for k, v in model.config.label2id.items()}
167     if list(sorted(label_name_to_id.keys())) == list(sorted(label_list)):
168         label_to_id = {i: int(label_name_to_id[label_list[i]]) for i in range(num_la
169     else:
170         logger.warning(
171             "Your model seems to have been trained with labels, but they don't match
172             f"model labels: {list(sorted(label_name_to_id.keys()))}, dataset labels:
173             "\nIgnoring the model labels as a result.",
174         )
175 elif data_args.task_name is None and not is_regression:
176     label_to_id = {v: i for i, v in enumerate(label_list)}

```

```
177
178     if label_to_id is not None:
179         model.config.label2id = label_to_id
180         model.config.id2label = {id: label for label, id in config.label2id.items()}
181     elif data_args.task_name is not None and not is_regression:
182         model.config.label2id = {l: i for i, l in enumerate(label_list)}
183         model.config.id2label = {id: label for label, id in config.label2id.items()}
184
185     if data_args.max_seq_length > tokenizer.model_max_length:
186         logger.warning(
187             f"The max_seq_length passed ({data_args.max_seq_length}) is larger than the
188             f'model ({tokenizer.model_max_length}). Using max_seq_length={tokenizer.mode
189         )
190     max_seq_length = min(data_args.max_seq_length, tokenizer.model_max_length)
191
192     def preprocess_function(examples):
193         # Tokenize the texts
194         args = (
195             (examples[sentence1_key],) if sentence2_key is None else (examples[sentence1
196         )
197         result = tokenizer(*args, padding=padding, max_length=max_seq_length, truncation
198
199         # Map labels to IDs (not necessary for GLUE tasks)
200         if label_to_id is not None and "label" in examples:
201             result["label"] = [(label_to_id[l] if l != -1 else -1) for l in examples["la
202         return result
203
204     with training_args.main_process_first(desc="dataset map pre-processing"):
205         raw_datasets = raw_datasets.map(
206             preprocess_function,
207             batched=True,
208             load_from_cache_file=not data_args.overwrite_cache,
209             desc="Running tokenizer on dataset",
210         )
211     if training_args.do_train:
212         if "train" not in raw_datasets:
213             raise ValueError("--do_train requires a train dataset")
214         train_dataset = raw_datasets["train"]
215         if data_args.max_train_samples is not None:
216             train_dataset = train_dataset.select(range(data_args.max_train_samples))
217
218     if training_args.do_eval:
219         if "validation" not in raw_datasets and "validation_matched" not in raw_datasets
220             raise ValueError("--do_eval requires a validation dataset")
221         eval_dataset = raw_datasets["validation_matched" if data_args.task_name == "mnli
222         if data_args.max_eval_samples is not None:
223             eval_dataset = eval_dataset.select(range(data_args.max_eval_samples))
224
225     if training_args.do_predict or data_args.task_name is not None or data_args.test_fil
226         if "test" not in raw_datasets and "test_matched" not in raw_datasets:
227             raise ValueError("--do_predict requires a test dataset")
```

```

228     predict_dataset = raw_datasets["test_matched" if data_args.task_name == "mnli" e
229     if data_args.max_predict_samples is not None:
230         predict_dataset = predict_dataset.select(range(data_args.max_predict_samples
231
232 # Log a few random samples from the training set:
233 if training_args.do_train:
234     for index in random.sample(range(len(train_dataset)), 3):
235         logger.info(f"Sample {index} of the training set: {train_dataset[index]}.")
236
237 # Get the metric function
238 if data_args.task_name is not None:
239     metric = load_metric("glue", data_args.task_name)
240 else:
241     metric = load_metric("accuracy")
242
243 # You can define your custom compute_metrics function. It takes an `EvalPrediction`
244 # predictions and label_ids field) and has to return a dictionary string to float.
245 def compute_metrics(p: EvalPrediction):
246     preds = p.predictions[0] if isinstance(p.predictions, tuple) else p.predictions
247     preds = np.squeeze(preds) if is_regression else np.argmax(preds, axis=1)
248     if data_args.task_name is not None:
249         result = metric.compute(predictions=preds, references=p.label_ids)
250         if len(result) > 1:
251             result["combined_score"] = np.mean(list(result.values())).item()
252         return result
253     elif is_regression:
254         return {"mse": ((preds - p.label_ids) ** 2).mean().item()}
255     else:
256         return {"accuracy": (preds == p.label_ids).astype(np.float32).mean().item()}
257
258 # Data collator will default to DataCollatorWithPadding when the tokenizer is passed
259 # we already did the padding.
260 if data_args.pad_to_max_length:
261     data_collator = default_data_collator
262 elif training_args.fp16:
263     data_collator = DataCollatorWithPadding(tokenizer, pad_to_multiple_of=8)
264 else:
265     data_collator = None
266
267 # Initialize our Trainer
268 trainer = Trainer(
269     model=model,
270     args=training_args,
271     train_dataset=train_dataset if training_args.do_train else None,
272     eval_dataset=eval_dataset if training_args.do_eval else None,
273     compute_metrics=compute_metrics,
274     tokenizer=tokenizer,
275     data_collator=data_collator,
276 )
277
278 # Training

```

```

279     if training_args.do_train:
280         checkpoint = None
281         if training_args.resume_from_checkpoint is not None:
282             checkpoint = training_args.resume_from_checkpoint
283         elif last_checkpoint is not None:
284             checkpoint = last_checkpoint
285         train_result = trainer.train(resume_from_checkpoint=checkpoint)
286         metrics = train_result.metrics
287         max_train_samples = (
288             data_args.max_train_samples if data_args.max_train_samples is not None else
289             )
290         metrics["train_samples"] = min(max_train_samples, len(train_dataset))
291
292         trainer.save_model() # Saves the tokenizer too for easy upload
293
294         trainer.log_metrics("train", metrics)
295         trainer.save_metrics("train", metrics)
296         trainer.save_state()
297
298     # Evaluation
299     if training_args.do_eval:
300         logger.info("*** Evaluate ***")
301
302         # Loop to handle MNLI double evaluation (matched, mis-matched)
303         tasks = [data_args.task_name]
304         eval_datasets = [eval_dataset]
305         if data_args.task_name == "mnli":
306             tasks.append("mnli-mm")
307             eval_datasets.append(raw_datasets["validation_mismatched"])
308
309         for eval_dataset, task in zip(eval_datasets, tasks):
310             metrics = trainer.evaluate(eval_dataset=eval_dataset)
311
312             max_eval_samples = (
313                 data_args.max_eval_samples if data_args.max_eval_samples is not None else
314                 )
315             metrics["eval_samples"] = min(max_eval_samples, len(eval_dataset))
316
317             trainer.log_metrics("eval", metrics)
318             trainer.save_metrics("eval", metrics)
319
320         kwargs = {"finetuned_from": model_args.model_name_or_path, "tasks": "text-classifica
321     if data_args.task_name is not None:
322         kwargs["language"] = "en"
323         kwargs["dataset_tags"] = "glue"
324         kwargs["dataset_args"] = data_args.task_name

```

▼ Run finetuning baselines

BERT is unstable and prone to degenerate performance on tasks with small training sets. The below cell fine-tunes BERT on tinySST (a small sentiment analysis dataset) using some default hyperparameters and also reports the mean and standard deviation of the dev set accuracy across 4 random seeds. Run the cell to obtain these baseline numbers, which should be around 50% average accuracy (it might take a couple of minutes to finish)

```

1 import timeit
2
3 start_time = timeit.default_timer()
4 task_name = "SST"
5 data_dir = f"./data/tiny{task_name}"
6 model_name_or_path = "bert-base-cased"
7 model_cache_dir = os.path.join(pretrained_models_dir, model_name_or_path)
8 data_cache_dir = f"./data_cache/finetuning/tiny{task_name}"
9
10 # Fine-tune and evaluate BERT with default hyperparameters using 4 random seeds
11 results = []
12 for seed in [1234, 2341, 3412, 4123]:
13     output_dir = f"./output/tiny{task_name}-{seed}"
14     config = dict(
15         seed=seed,
16         model_name_or_path=model_name_or_path,
17         train_file="./data/tinySST/train.csv",
18         validation_file="./data/tinySST/dev.csv",
19         task_type="text_classification",
20         do_train=True,
21         do_eval=True,
22         do_lower_case=True,
23         data_dir=data_dir,
24         max_seq_length=128,
25         per_device_train_batch_size=32,
26         learning_rate=2e-5,
27         num_train_epochs=3.0,
28         model_cache_dir=model_cache_dir,
29         data_cache_dir=data_cache_dir,
30         output_dir=output_dir,
31         overwrite_output_dir=True,
32         log_level='warning'
33     )
34
35     result = do_target_task_finetuning(config)
36     results.append(result["eval_accuracy"])
37
38 results = np.array(results)
39 mean = np.mean(results)
40 std = np.std(results)
41
42 print(f"Accuracy on TinySST dev set: {mean} +/- {std}")
43 elapsed_time = timeit.default_timer() - start_time
44 print(f"Time elapsed: {elapsed_time} seconds")

```



```

WARNING:__main__:Process rank: -1, device: cuda:0, n_gpu: 1distributed training: False,
WARNING:datasets.builder:Using custom data configuration default-ae0aa6a33dbb7123
Downloading and preparing dataset csv/default to /root/.cache/huggingface/datasets/csv/
Downloading data files: 100% 2/2 [00:00<00:00, 64.36it/s]

Extracting data files: 100% 2/2 [00:00<00:00, 52.54it/s]
Dataset csv downloaded and prepared to /root/.cache/huggingface/datasets/csv/default-ae
100% 2/2 [00:00<00:00, 61.43it/s]

Downloading: 100% 570/570 [00:00<00:00, 15.6kB/s]
Downloading: 100% 29.0/29.0 [00:00<00:00, 1.12kB/s]
Downloading: 100% 208k/208k [00:00<00:00, 245kB/s]
Downloading: 100% 426k/426k [00:01<00:00, 449kB/s]
Downloading: 100% 416M/416M [00:25<00:00, 16.6MB/s]

[WARNING|modeling_utils.py:1693] 2022-12-03 16:15:18,445 >> Some weights of the model c
- This IS expected if you are initializing BertForSequenceClassification from the check
- This IS NOT expected if you are initializing BertForSequenceClassification from the c
[WARNING|modeling_utils.py:1704] 2022-12-03 16:15:18,447 >> Some weights of BertForSequ
You should probably TRAIN this model on a down-stream task to be able to use it for pre
Running tokenizer on dataset: 1/1 [00:00<00:00,
100% 25.81ba/s]

Running tokenizer on dataset: 1/1 [00:00<00:00,
100% 6.85ba/s]

Downloading builder script: 3.19k/? [00:00<00:00, 66.1kB/s]

```

 [3/3 00:00, Epoch 3/3]

Step Training Loss

***** train metrics *****

```

epoch = 3.0
total_flos = 3675GF
train_loss = 0.6892
train_runtime = 0:00:01.21
train_samples = 20
train_samples_per_second = 49.21
train_steps_per_second = 2.46

```

 [109/109 00:06]

```

WARNING:__main__:Process rank: -1, device: cuda:0, n_gpu: 1distributed training: False,
***** eval metrics *****

```

```

epoch = 3.0
eval_accuracy = 0.4977
eval_loss = 0.6911
eval_runtime = 0:00:06.20
eval_samples = 872
eval_samples_per_second = 140.541
eval_steps_per_second = 17.568

```

```

WARNING:datasets.builder:Using custom data configuration default-ae0aa6a33dbb7123

```

```

WARNING:datasets.builder:Reusing dataset csv (/root/.cache/huggingface/datasets/csv/def

```

100%

2/2 [00:00<00:00, 50.16it/s]

[WARNING|modeling_utils.py:1693] 2022-12-03 16:15:41,168 >> Some weights of the model c
 - This IS expected if you are initializing BertForSequenceClassification from the check
 - This IS NOT expected if you are initializing BertForSequenceClassification from the c
 [WARNING|modeling_utils.py:1704] 2022-12-03 16:15:41,169 >> Some weights of BertForSequ
 You should probably TRAIN this model on a down-stream task to be able to use it for pre

Running tokenizer on dataset:

1/1 [00:00<00:00,

100%

27.43ba/s]

Running tokenizer on dataset:

1/1 [00:00<00:00,

100%

5.23ba/s]

 [3/3 00:00, Epoch 3/3]

Step Training Loss

***** train metrics *****

epoch	=	3.0
total_flos	=	3675GF
train_loss	=	0.7235
train_runtime	=	0:00:01.23
train_samples	=	20
train_samples_per_second	=	48.69
train_steps_per_second	=	2.434

 [109/109 00:06]

WARNING:__main__:Process rank: -1, device: cuda:0, n_gpu: 1distributed training: False,

***** eval metrics *****

epoch	=	3.0
eval_accuracy	=	0.5
eval_loss	=	0.6947
eval_runtime	=	0:00:06.26
eval_samples	=	872
eval_samples_per_second	=	139.199
eval_steps_per_second	=	17.4

WARNING:datasets.builder:Using custom data configuration default-ae0aa6a33dbb7123

WARNING:datasets.builder:Reusing dataset csv (/root/.cache/huggingface/datasets/csv/def

100%

2/2 [00:00<00:00, 67.32it/s]

[WARNING|modeling_utils.py:1693] 2022-12-03 16:16:03,647 >> Some weights of the model c
 - This IS expected if you are initializing BertForSequenceClassification from the check
 - This IS NOT expected if you are initializing BertForSequenceClassification from the c
 [WARNING|modeling_utils.py:1704] 2022-12-03 16:16:03,654 >> Some weights of BertForSequ
 You should probably TRAIN this model on a down-stream task to be able to use it for pre

Running tokenizer on dataset:

1/1 [00:00<00:00,

100%


25.24ba/s]

Running tokenizer on dataset:

1/1 [00:00<00:00,

100%

5.85ba/s]

 [3/3 00:00, Epoch 3/3]

Step Training Loss

***** train metrics *****

epoch	=	3.0
-------	---	-----

```

total_flos          = 3675GF
train_loss          = 0.6883
train_runtime       = 0:00:01.19
train_samples       = 20
train_samples_per_second = 50.061
train_steps_per_second   = 2.503

```

[109/109 00:06]

WARNING:__main__:Process rank: -1, device: cuda:0, n_gpu: 1distributed training: False,
 ***** eval metrics *****

```

epoch              = 3.0
eval_accuracy      = 0.6055
eval_loss          = 0.6732
eval_runtime       = 0:00:06.32
eval_samples       = 872
eval_samples_per_second = 137.759
eval_steps_per_second   = 17.22

```

WARNING:datasets.builder:Using custom data configuration default-ae0aa6a33dbb7123

WARNING:datasets.builder:Reusing dataset csv (/root/.cache/huggingface/datasets/csv/default-ae0aa6a33dbb7123)
 100% 2/2 [00:00<00:00, 67.06it/s]

[WARNING|modeling_utils.py:1693] 2022-12-03 16:16:26,281 >> Some weights of the model c
 - This IS expected if you are initializing BertForSequenceClassification from the check
 - This IS NOT expected if you are initializing BertForSequenceClassification from the c
 [WARNING|modeling_utils.py:1704] 2022-12-03 16:16:26,283 >> Some weights of BertForSequ
 You should probably TRAIN this model on a down-stream task to be able to use it for pre

Running tokenizer on dataset: 1/1 [00:00<00:00,

100% 31.20ba/s]

Running tokenizer on dataset: 1/1 [00:00<00:00,

100% 5.58ba/s]

[3/3 00:00, Epoch 3/3]

Step Training Loss

***** train metrics *****

```

epoch              = 3.0
total_flos          = 3675GF
train_loss          = 0.6595
train_runtime       = 0:00:01.23
train_samples       = 20
train_samples_per_second = 48.602

```

▼ Run translate demo

Now run the following cell to load Google Translate's model and run it on a toy example. You will use Google Translate to augment your TinySST dataset via backtranslation, which involves translating an example to another language (or languages) and then eventually translating it back to English. This process injects syntactic and lexical variation into the input which can help the model learn.

```

1 import googletrans
2 # Run print(googletrans.LANGUAGES) to see available languages
3 from googletrans import Translator
4 translator = Translator()
5
6 # translate from English to French
7 output = translator.translate("I love natural language processing", src='en', dest='fa')
8 output.text

```

'من عاشق پردازش زبان طبیعی هستم'

▼ Question 3.1 (20 points)

Complete the following cell to paraphrase the training data of `tinySST` using backtranslation. We have intentionally left this problem open-ended: feel free to use as many pivot languages as you like, and also write any postprocessing code you think might help. The cell after this one will fine-tune BERT on the augmented training data, so you can use its output to validate your backtranslation strategy. To obtain full points, the model fine-tuned on your augmented data must achieve a higher average accuracy (averaged across random seeds) than the model without any augmentation, trained with the same hyperparameters.

Write your code in the indicated section in the below cell.

```

1 task_name = "SST"
2 data_dir = f"./data/tiny{task_name}"
3 task_processor = glue_processors[f"{task_name.lower()}-2"]()
4 train_examples = task_processor.get_train_examples(data_dir)
5
6 train_examples_augmented = []
7
8 ### (incomplete) list of languages you can use
9 languages = [
10     'en', # english
11     'cs', # czech
12     'de', # german
13     'es', # spanish
14     'fi', # finnish
15     'fr', # french
16     'hi', # hindi
17     'it', # italian
18     'ja', # japanese
19     'pt', # portuguese
20     'ru', # russian
21     'vi', # vietnamese
22     'zh-cn', # chinese
23     'fa', # Persian
24 ]

```

```
25 PIVOT_LANGUAGES = ['en', 'cs', 'de', 'es', 'fi', 'fr', 'hi', 'it', 'ja', 'pt', 'ru', 'vi']
26
27 # generate some augmented examples for each training example
28 for example in tqdm(train_examples):
29     train_examples_augmented.append(example) # always include the original example
30     print(example)
31     # WRITE YOUR CODE HERE!
32     for target_language in PIVOT_LANGUAGES:
33         pivot = translator.translate(example.text_a, src='en', dest=target_language).text
34         paraphrase = translator.translate(pivot, src=target_language, dest='en').text
35         # the below line adds a single new augmented example to the dataset.
36         # note that the guid should be a unique ID for this example, so you'll want to v
37         # depending on how you generate your paraphrases
38         train_examples_augmented.append(InputExample(guid=f"{example.guid}-aug-{target_l
39                                                         text_a=paraphrase,
40                                                         text_b=None,
41                                                         label=example.label))
42 output_dir = f"./data/tiny{task_name}-bt"
43 if not os.path.exists(output_dir):
44     os.makedirs(output_dir)
45
46 with open(os.path.join(output_dir, "train.tsv"), "w") as writer:
47     writer.write("sentence\tlabel\n")
48     for example in train_examples_augmented:
49         writer.write(f"{example.text_a}\t{example.label}\n")
50 tsv_to_csv(os.path.join(output_dir, "train.tsv"), os.path.join(output_dir, "train.csv"))
51
52 # Copy the original tinySST's dev set to the new directory
53 import shutil
54 shutil.copyfile(f"{data_dir}/dev.csv", f"{output_dir}/dev.csv")
```

```

a%l      | a/2a [aa·aa/2 2it/εlTnputExample(guid='train-1' text a='its unerring
1 # Examples
2 for i in range(len(PIVOT_LANGUAGES) + 1):
3     print(train_examples_augmented[i].guid.ljust(20), train_examples_augmented[i].text_a

train-1          its unerring respect for them
train-1-aug-en    its unerring respect for them
train-1-aug-cs    his unfailing respect for them
train-1-aug-de    his unfailing respect for her
train-1-aug-es    his unfailing respect for them
train-1-aug-fi    its unmistakable respect for them
train-1-aug-fr    his unfailing respect for them
train-1-aug-hi    it's an absolute honor for him
train-1-aug-it    his unfailing respect for them
train-1-aug-ja    unwavering respect for them
train-1-aug-pt    his unfailing respect for them
train-1-aug-ru    his unmistakable respect for them
train-1-aug-vi    its unwavering respect for them
train-1-aug-zh-cn its respect for them
train-1-aug-fa    Its unparalleled respect for them

' /data/tinySST-bt/dev.csv'

```

The below cell fine-tunes BERT bert-base-cased with the combined training data (real + synthetic training examples) and then evaluates the resulting model on tinySST's dev set. Note that it uses the default fine-tuning hyperparameters, not the improved ones that you found earlier. You should observe a significantly higher accuracy than 50% when you run this cell on the augmented data (our reference implementation reaches 64%). **Do NOT modify any code in this cell!**

```

1 import timeit
2
3 start_time = timeit.default_timer()
4 task_name = "SST"
5 data_dir = f"./data/tiny{task_name}-bt"
6 model_name_or_path = "bert-base-cased"
7 model_cache_dir = os.path.join(pretrained_models_dir, model_name_or_path)
8 data_cache_dir = f"./data_cache/finetuning/tiny{task_name}-bt/"
9 output_dir = model_cache_dir
10
11 # Fine-tune and evaluate BERT with default hyperparameters using 4 random seeds
12 results = []
13 for seed in [1234, 2341, 3412, 4123]:
14     output_dir = f"./output/tiny{task_name}-{seed}"
15     config = dict(
16         seed=seed,
17         model_name_or_path=model_name_or_path,
18         train_file="./data/tinySST-bt/train.csv",
19         validation_file="./data/tinySST-bt/dev.csv",
20         task_type="text_classification",
21         do_train=True,
22         do_eval=True,

```

```
23     do_lower_case=True,
24     data_dir=data_dir,
25     max_seq_length=128,
26     per_device_train_batch_size=32,
27     learning_rate=2e-5,
28     num_train_epochs=3.0,
29     model_cache_dir=model_cache_dir,
30     data_cache_dir=data_cache_dir,
31     output_dir=output_dir,
32     overwrite_output_dir=True,
33     log_level='warning'
34 )
35
36 result = do_target_task_finetuning(config)
37 results.append(result["eval_accuracy"])
38
39 results = np.array(results)
40 mean = np.mean(results)
41 std = np.std(results)
42
43 print(f"Accuracy on TinySST dev set: {mean} +/- {std}")
44 elapsed_time = timeit.default_timer() - start_time
45 print(f"Time elapsed: {elapsed_time} seconds")
```

```

WARNING:__main__:Process rank: -1, device: cuda:0, n_gpu: 1distributed training: False,
WARNING:datasets.builder:Using custom data configuration default-8f825126d2882a66
Downloading and preparing dataset csv/default to /root/.cache/huggingface/datasets/csv/
Downloading data files: 100%                2/2 [00:00<00:00, 70.79it/s]

Extracting data files: 100%                2/2 [00:00<00:00, 38.53it/s]

Dataset csv downloaded and prepared to /root/.cache/huggingface/datasets/csv/default-8f
100%                2/2 [00:00<00:00, 66.04it/s]

[WARNING|modeling_utils.py:1693] 2022-12-03 16:46:24,053 >> Some weights of the model c
- This IS expected if you are initializing BertForSequenceClassification from the check
- This IS NOT expected if you are initializing BertForSequenceClassification from the c
[WARNING|modeling_utils.py:1704] 2022-12-03 16:46:24,055 >> Some weights of BertForSequ
You should probably TRAIN this model on a down-stream task to be able to use it for pre

Running tokenizer on dataset:                1/1 [00:00<00:00,
100%                14.56ba/s]

Running tokenizer on dataset:                1/1 [00:00<00:00,
100%                5.86ba/s]

/usr/local/lib/python3.8/dist-packages/transformers/optimization.py:306: FutureWarning:
warnings.warn(
[30/30 00:16, Epoch 3/3]

```

Step Training Loss

***** train metrics *****

```

epoch                =          3.0
total_flos           =      55134GF
train_loss           =       0.4333
train_runtime        = 0:00:16.73
train_samples        =       300
train_samples_per_second =      53.795
train_steps_per_second  =       1.793

```

[109/109 00:06]

```

WARNING:__main__:Process rank: -1, device: cuda:0, n_gpu: 1distributed training: False,
***** eval metrics *****

```

```

epoch                =          3.0
eval_accuracy        =       0.6927
eval_loss            =       0.5946
eval_runtime         = 0:00:06.39
eval_samples         =       872
eval_samples_per_second =     136.278
eval_steps_per_second  =      17.035

```

```

WARNING:datasets.builder:Using custom data configuration default-8f825126d2882a66
WARNING:datasets.builder:Reusing dataset csv (/root/.cache/huggingface/datasets/csv/def
100%                2/2 [00:00<00:00, 68.05it/s]

```

```

[WARNING|modeling_utils.py:1693] 2022-12-03 16:47:02,443 >> Some weights of the model c
- This IS expected if you are initializing BertForSequenceClassification from the check
- This IS NOT expected if you are initializing BertForSequenceClassification from the c
[WARNING|modeling_utils.py:1704] 2022-12-03 16:47:02,445 >> Some weights of BertForSequ
You should probably TRAIN this model on a down-stream task to be able to use it for pre

Running tokenizer on dataset:                1/1 [00:00<00:00,

```


100%


8.41ba/s]

Running tokenizer on dataset:

1/1 [00:00<00:00,

100%

8.14ba/s]

 [30/30 00:16, Epoch 3/3]**Step Training Loss**

***** train metrics *****

epoch	=	3.0
total_flos	=	55134GF
train_loss	=	0.5299
train_runtime	=	0:00:17.15
train_samples	=	300
train_samples_per_second	=	52.469
train_steps_per_second	=	1.749

 [109/109 00:06]

WARNING:__main__:Process rank: -1, device: cuda:0, n_gpu: 1distributed training: False,

***** eval metrics *****

epoch	=	3.0
eval_accuracy	=	0.6938
eval_loss	=	0.6069
eval_runtime	=	0:00:06.63
eval_samples	=	872
eval_samples_per_second	=	131.479
eval_steps_per_second	=	16.435

WARNING:datasets.builder:Using custom data configuration default-8f825126d2882a66

WARNING:datasets.builder:Reusing dataset csv (/root/.cache/huggingface/datasets/csv/def

100%

2/2 [00:00<00:00, 29.35it/s]

[WARNING|modeling_utils.py:1693] 2022-12-03 16:47:42,654 >> Some weights of the model c
 - This IS expected if you are initializing BertForSequenceClassification from the check
 - This IS NOT expected if you are initializing BertForSequenceClassification from the c
 [WARNING|modeling_utils.py:1704] 2022-12-03 16:47:42,657 >> Some weights of BertForSequ
 You should probably TRAIN this model on a down-stream task to be able to use it for pre

Running tokenizer on dataset:

1/1 [00:00<00:00,

100%

6.17ba/s]

Running tokenizer on dataset:

1/1 [00:00<00:00,

100%

5.98ba/s]

 [30/30 00:18, Epoch 3/3]**Step Training Loss**

***** train metrics *****

epoch	=	3.0
total_flos	=	55134GF
train_loss	=	0.442
train_runtime	=	0:00:18.86
train_samples	=	300
train_samples_per_second	=	47.704
train_steps_per_second	=	1.59

 [109/109 00:07]

WARNING:__main__:Process rank: -1, device: cuda:0, n_gpu: 1distributed training: False,

***** eval metrics *****

```

epoch                =      3.0
eval_accuracy        =      0.6709
eval_loss            =      0.5946
eval_runtime         = 0.00.07 19

```

▼ Question 3.2 (5 points)

Briefly explain your backtranslation strategy here. Why do you think it resulted in an improvement?

100%

2/2 100.00-00.00 20.40s/1

Write your answer here! Please keep it brief (i.e., 2-3 sentences).

First I observed the augmented examples for each pivot language:

train-1	its unerring respect for them
train-1-aug-en	its unerring respect for them
train-1-aug-cs	his unfailing respect for them
train-1-aug-de	his unfailing respect for her
train-1-aug-es	his unfailing respect for them
train-1-aug-fi	its unmistakable respect for them
train-1-aug-fr	his unfailing respect for them
train-1-aug-hi	it's an absolute honor for him
train-1-aug-it	his unfailing respect for them
train-1-aug-ja	unwavering respect for them
train-1-aug-pt	his unfailing respect for them
train-1-aug-ru	his unmistakable respect for them
train-1-aug-vi	its unwavering respect for them
train-1-aug-zh-cn	its respect for them
train-1-aug-fa	Its unparalleled respect for them

As seen above, each backtranslation has created a new instance for us with new words and even in some cases such as hindi (hi) a new sentence structure. These new variated samples will help our model to better generalize and learn the task instead of overfitting to specific words. (The sentiment prediction task is good for this appraoch as long as the backtranslated versions do not change the sentiment which is rare to happen.)

Accuracy on TinySST dev set: 0.6751720309257507 +/- 0.020520400137613005

[Colab paid products](#) - [Cancel contracts here](#)

✓ 2m 43s completed at 8:18 PM ● ✕