

Natural Language Processing

HomeWork 4

Ilia Hashemi Rad

99102456

AmirMohammad Fakhimi

99170531

AmirMahdi Namjoo

402211467

Introduction

Explanations:

In this project, we have implemented a Transformer model that tries to identify and neutralize gender bias in texts. To achieve this, we first gathered a dataset consisting of Persian tweets with gender bias and some ChatGPT responses with gender bias. We then investigated them manually. We then used LLMs to neutralize gender bias in these texts. We then fine-tuned an LLM on these data and evaluated the results in each part.

We should also mention that we used "Identifying and Neutralizing Gender Bias from Text" from Stanford University to get some ideas for the implementation. Notably, the article is based on the English language, and we did not have access to its dataset or code.

It is noteworthy that a dataset called "EXIST" exists for gender bias in the English and Spanish language, but we could not get it as it needs filling out a form to get authorized for access.

Collecting and Preparing the Dataset

Explanations:

As there is no Persian dataset for the gender bias problem, we gathered the dataset by scrapping Twitter using "selenium-twitter-scraper". The amount of data collected from Twitter was not enough to get appropriate results, so we asked ChatGPT to generate another set of data with and without gender bias to increase the size of our dataset significantly.

After getting the data, we processed them to prepare them for use in the models in the next section of this assignment.

Note: Some parts of the explanations of code from here to the end of the document are written by partially ChatGPT to help us write clear and unambiguous descriptions.

```
1 try :
2     os.makedirs('data/tweets')
3 except FileExistsError:
4     print('Directory already exists')
5 os.chdir('selenium-twitter-scraper')
6
7 base_command = f'python scraper --user=@{config["TWITTER_USERNAME"]} --password={config["TWITTER_PASSWORD"]} --top '
8
9 os.system(base_command + '--tweets 150 --query "***'')
10 os.system(base_command + '--tweets 150 --query "***'')
11 os.system(base_command + '--tweets 150 --query "***'')
12 os.system(base_command + '--tweets 150 --query "***'')
13 os.system(base_command + '--tweets 150 --query "***'')
14
15 os.system('cp -r tweets ../data')
16 os.chdir('..')
17 directory = 'data/tweets'
18 tweets_files = os.listdir(directory)
19 tweets_files = [tweets_file for tweets_file in tweets_files if tweets_file.startswith('2024')]
20 tweets_files = sorted(tweets_files)
21 tweets = []
22
23 for tweets_file in tweets_files:
24     current_file = pd.read_csv(f'{directory}/{tweets_file}')
25     current_file = current_file['Content'].tolist()
26     tweets += current_file
27
28 len(tweets), tweets
```

Explanations:

The main part of scrapping Twitter is using the Python scraper command with different queries. Note that due to L^AT_EX limitations in code, we could not put Persian text in the code section above, but the queries are Persian equivalents of "Feminism", "Only Boy", "Only Girl", "Only Man", "Only Woman". After that, we get only those that are from 2024 and store them.

```
1 directory = 'data/chatgpt'
2
3 biased_files = [
4     'Bias1.txt',
5     'Bias2.txt',
6     'Bias3.txt',
7     'Bias4.txt',
8     'Bias5.txt',
9 ]
10
11 unbiased_files = [
12     'NoBias1.txt',
13     'NoBias2.txt',
14     'NoBias3.txt',
15     'NoBias4.txt',
16     # 'NoBias5.txt',
17     'NoBias6.txt',
18 ]
19
20 biased_chatgpt = []
21 unbiased_chatgpt = []
22
23 for file in biased_files:
24     with open(f'{directory}/{file}', 'r', encoding='utf-8') as file:
25         current_data = file.read().split('\n\n')
26         biased_chatgpt += current_data
27
28 for file in unbiased_files:
29     with open(f'{directory}/{file}', 'r', encoding='utf-8') as file:
30         current_data = file.read().split('\n\n')
31         unbiased_chatgpt += current_data
32
33 len(biased_chatgpt), len(unbiased_chatgpt)
```

Explanations:

The script reads biased and unbiased data files:

- Defines lists of file names for biased and unbiased data.
- Reads and splits the content of each file into lists.

```
1 data = []
2
3 with open('data/tweets/tweets_all.txt', 'r', encoding='utf-8') as file:
4     labeled_tweets = file.read().split('\n\n')
5
6 len(labeled_tweets)
7
8 for i in range(0, len(labeled_tweets), 2):
9     tweet, label = labeled_tweets[i], labeled_tweets[i + 1]
10    data.append({
11        "id": len(data),
12        "dialogue": tweet,
13        "summary": label,
14    })
15
16 for biased_chatgpt_data in biased_chatgpt:
```

```

17 prompt, label = biased_chatgpt_data.split('\n')
18 prompt = prompt.split(':')[1]
19 label = label.split(':')[1]
20 data.append({
21     "id": len(data),
22     "dialogue": prompt,
23     "summary": label,
24 })
25
26 for unbiased_chatgpt_data in unbiased_chatgpt:
27     prompt = unbiased_chatgpt_data
28     prompt = prompt.split(':')[1]
29     data.append({
30         "id": len(data),
31         "dialogue": prompt,
32         "summary": prompt,
33     })

```

Explanations:

Finally, the script combines labeled tweet data and ChatGPT data into a single list of dictionaries:

- Reads labeled tweets from a file.
- Appends tweet and label pairs to the data list.
- Processes and appends biased and unbiased ChatGPT data.

```

1 with open('stopwords/chars.txt', 'r', encoding='utf-8') as file:
2     chars_stop_words = ''.join(file.read().splitlines())
3
4     chars_stop_words = chars_stop_words.replace('[', '\[')
5     chars_stop_words = chars_stop_words.replace(']', '\]')
6     chars_stop_words += '\t\r'
7     chars_pattern = re.compile(f'[{chars_stop_words}]')
8     chars_pattern

```

Explanations:

The script reads stopwords from a file and prepares a regex pattern to remove them:

- Reads stopwords from `chars.txt` file.
- Modifies stopwords string to escape special characters.
- Compiles a regex pattern to match these stopwords.

```

1 # https://stackoverflow.com/questions/33404752/removing-emojis-from-a-string-in-python
2 emojis_pattern = re.compile("[
3     u"\U0001F600-\U0001F64F" # emoticons
4     u"\U0001F300-\U0001F5FF" # symbols & pictographs
5     u"\U0001F680-\U0001F6FF" # transport & map symbols
6     u"\U0001F1E0-\U0001F1FF" # flags (iOS)
7     "]+")
8 emojis_pattern

```

Explanations:

The script prepares a regex pattern to remove emojis, referencing a solution from Stack Overflow.

```

1 normalizer = Normalizer()

```

Explanations:

The script initializes a `Normalizer` from the `hazm` library for text normalization.

```
1 def remove_duplicates(data):
2     data_df: pd.DataFrame = pd.DataFrame(data)
3     data_df = data_df.drop_duplicates(subset=['dialogue'])
4
5     data = []
6     for index, row in data_df.iterrows():
7         data.append({
8             "id": row['id'],
9             "dialogue": row['dialogue'],
10            "summary": row['summary'],
11        })
12
13     return data
```

Explanations:

The `remove_duplicates` function removes duplicate entries from the data based on the `dialogue` field.

```
1 def preprocess_text(text):
2     global chars_pattern, emojis_pattern, normalizer
3
4     text = emojis_pattern.sub(r' ', text)
5     text = chars_pattern.sub(r' ', text)
6     text = normalizer.normalize(text)
7
8     return text
```

Explanations:

The `preprocess_text` function applies regex patterns to remove emojis and stopwords and normalizes the text.

```
1 len_before = len(data)
2 data = remove_duplicates(data)
3 len_after = len(data)
4
5 len_before, len_after
```

Explanations:

The script removes duplicates from the data and prints the length of the data before and after the removal.

```
1 preprocessed_data = []
2
3 progress_bar = tqdm(data)
4 for datum in progress_bar:
5     preprocessed_data.append({
6         "id": datum["id"],
7         "dialogue": preprocess_text(datum["dialogue"]),
8         "summary": preprocess_text(datum["summary"]),
9     })
```

Explanations:

Finally, the script preprocesses the text data using the `preprocess_text` function and stores the results in `preprocessed_data`.

```

1 import py7zr
2 import json
3 from sklearn.model_selection import train_test_split
4 import random
5
6 NUMBER_TO_SHUFFLE = 4
7 random.seed(42)
8
9 for _ in range(NUMBER_TO_SHUFFLE):
10     random.shuffle(preprocessed_data)
11     print(preprocessed_data[0])
12
13 train_data, test_data = train_test_split(preprocessed_data, test_size=0.1, random_state=42)
14 len(train_data), len(test_data)
15
16 train_data = sorted(train_data, key=lambda x: x['id'])
17 test_data = sorted(test_data, key=lambda x: x['id'])

```

Explanations:

This segment of the script performs the following operations:

- `import py7zr`: Imports the `py7zr` library for handling 7z archive files.
- `import json`: Imports the `json` library for handling JSON data.
- `from sklearn.model_selection import train_test_split`: Imports the `train_test_split` function from `sklearn` for splitting data into training and testing sets.
- `import random`: Imports the `random` library for random operations.

The script then sets up the shuffling and splitting of preprocessed data:

- `NUMBER_TO_SHUFFLE = 4`: Defines the number of times the data will be shuffled.
- `random.seed(42)`: Sets the random seed for reproducibility.
- A `for` loop shuffles the `preprocessed_data` list `NUMBER_TO_SHUFFLE` times and prints the first element after each shuffle.
- `train_data, test_data = train_test_split(preprocessed_data, test_size=0.1, random_state=42)`: Splits the data into training and testing sets, with 10% of the data reserved for testing.
- `train_data = sorted(train_data, key=lambda x: x['id'])`: Sorts the training data by the `id` field.
- `test_data = sorted(test_data, key=lambda x: x['id'])`: Sorts the testing data by the `id` field.

```

1 directory = 'gender_neutralize/data'
2
3 os.makedirs(directory, exist_ok=True)
4
5 with open(f'{directory}/train.json', 'w', encoding='utf-8') as file:
6     json.dump(train_data, file, ensure_ascii=False, indent=2)
7
8 with open(f'{directory}/test.json', 'w', encoding='utf-8') as file:
9     json.dump(test_data, file, ensure_ascii=False, indent=2)
10
11 os.chdir(directory)
12
13 with py7zr.SevenZipFile('gender_neutralize.7z', 'w') as archive:
14     archive.write('train.json')
15     archive.write('test.json')

```

```
16 os.chdir(' ../../..')
17
```

Explanations:

This part of the script saves and archives the processed data:

- `directory = 'gender_neutralize/data'`: Specifies the directory to store the data.
- `os.makedirs(directory, exist_ok=True)`: Creates the directory if it does not already exist.
- The script then opens two files, `train.json` and `test.json`, in write mode and dumps the training and testing data into these files as JSON objects.
- `os.chdir(directory)`: Changes the current working directory to the specified directory.
- `with py7zr.SevenZipFile('gender_neutralize.7z', 'w') as archive`: Creates a new 7z archive file named `gender_neutralize.7z` in write mode.
- `archive.write('train.json')`: Adds the `train.json` file to the archive.
- `archive.write('test.json')`: Adds the `test.json` file to the archive.
- `os.chdir(' ../../..')`: Changes the current working directory back to the original directory.

Fine Tuning

In this section, we will explain the code for using LLaMa 3 and fine-tuning it for the neutralization of gender bias.

We used LLaMa 3 with 8Bilion Parameters. There is also a 70 Billion Parameter version, but we removed it from our options as it needs more resources. We studied some of its tutorials that fine-tuned LLaMa 3 on a summarization task, and we used the ideas and code for our task, although with some changes. We used the summarizer structure as its input-output format, much like our task format; the summarizer task gets an input text and outputs it as a summary, and in our task, we have an input text and need the gender-bias-neutralized text as the output.

The code fine-tunes our LLaMa3 using our dataset gathered in the previous parts using LoRA. An important note here is that this task needs more than 16 GB of GPU Memory, and we did not have access to this in Colab. Fortunately, the Kaggle enabled us to run our model on two T4 accelerators in parallel, solving our problem of limited resources on Colab. Also, we changed the context-length parameter from 4096 to 1024 to mitigate the lack of resources.

Also, after fine-tuning, we saved the model, emptied the GPU memory, and loaded the model again. This helped us mitigate some problems regarding insufficient memory.

Now, let's get to the code.

Step 0

The following code installs necessary Python packages and logs into the Hugging Face Hub.

```
! pip install --upgrade pip
! pip install evaluate
! pip install bert_score
! pip install hazm
! pip install llama-recipes ipywidgets

import huggingface_hub
huggingface_hub.login()
```

Step 1

```
1 import torch
2 from transformers import LlamaForCausalLM, AutoTokenizer
3 from llama_recipes.configs import train_config as TRAIN_CONFIG
4
5 train_config = TRAIN_CONFIG()
6 train_config.model_name = "meta-llama/Meta-Llama-3-8B"
7 train_config.num_epochs = 3
8 train_config.run_validation = False
9 train_config.gradient_accumulation_steps = 4
10 train_config.batch_size_training = 1
11 train_config.lr = 2e-4
12 train_config.use_fast_kernels = True
13 train_config.use_fp16 = True
14 train_config.context_length = 1024 if torch.cuda.get_device_properties(0).total_memory < 16e9
15     else 2048 # T4 16GB or A10 24GB
16 train_config.batching_strategy = "packing"
17 train_config.output_dir = "meta-llama-sexbias"
18
19 from transformers import BitsAndBytesConfig
20 config = BitsAndBytesConfig(
21     load_in_8bit=True,
22 )
23
24 model = LlamaForCausalLM.from_pretrained(
25     train_config.model_name,
26     device_map="auto",
27     quantization_config=config,
28     use_cache=False,
29     attn_implementation="sdpa" if train_config.use_fast_kernels else None,
30     torch_dtype=torch.float16,
31 )
32
33 tokenizer = AutoTokenizer.from_pretrained(train_config.model_name)
34 tokenizer.pad_token = tokenizer.eos_token
```

Explanations:

First, import the necessary libraries. `torch` is used for tensor computations, `LlamaForCausalLM` and `AutoTokenizer` are used for model and tokenizer management from the `transformers` library, and `train_config` is imported from `llama_recipes.configs`.

```
1 import torch
2 from transformers import LlamaForCausalLM, AutoTokenizer
3 from llama_recipes.configs import train_config as TRAIN_CONFIG
```

Explanations:

Next, create a training configuration object and set various training parameters. These include the model name, number of epochs, validation settings, gradient accumulation steps, batch size for training, learning rate, kernel usage, floating point precision, context length based on available GPU memory, batching strategy, and output directory.

```
1 train_config = TRAIN_CONFIG()
2 train_config.model_name = "meta-llama/Meta-Llama-3-8B"
3 train_config.num_epochs = 3
4 train_config.run_validation = False
5 train_config.gradient_accumulation_steps = 4
6 train_config.batch_size_training = 1
7 train_config.lr = 2e-4
8 train_config.use_fast_kernels = True
9 train_config.use_fp16 = True
10 train_config.context_length = 1024 if torch.cuda.get_device_properties(0).total_memory
11     < 16e9 else 2048
```



```

1 train_config.batching_strategy = "packing"
2 train_config.output_dir = "meta-llama-sexbias"

```

Explanations:

Import the `BitsAndBytesConfig` class from the `transformers` library and create a configuration object for loading the model in 8-bit precision.

```

1 from transformers import BitsAndBytesConfig
2 config = BitsAndBytesConfig(
3     load_in_8bit=True,
4 )

```

Explanations:

Load the LLaMA model for causal language modeling using the specified configuration. The model is loaded with the training configuration parameters, including device mapping, quantization settings, cache usage, attention implementation, and data type.

```

1 model = LlamaForCausalLM.from_pretrained(
2     train_config.model_name,
3     device_map="auto",
4     quantization_config=config,
5     use_cache=False,
6     attn_implementation="sdpa" if train_config.use_fast_kernels else None,
7     torch_dtype=torch.float16,
8 )

```

Explanations:

Load the tokenizer associated with the LLaMA model and set the padding token to the end-of-sequence token.

```

1 tokenizer = AutoTokenizer.from_pretrained(train_config.model_name)
2 tokenizer.pad_token = tokenizer.eos_token

```

Step 2

The following code prepares a dataset, generates predictions using a model, computes evaluation metrics, and performs a sample evaluation.

As the Prompt template is a Persian text that could not be adequately rendered in the \LaTeX listing environment, we have included its picture below and put `***` in the code on this document. It can be appropriately viewed in the Jupyter Notebook. The best way is to ask, "Remove this text gender bias and make it without gender bias," then add the text and put "text without gender bias:" with eight underlines in Persian.

It is also noteworthy that we used the Hazm tokenizer. It produces better results than not using a tokenizer at all. The base code of Facebook examples uses the nltk tokenizer (for English).

For the Bert-Score calculation, we used Microsoft `deberta-v2-xxlarge-mnli`.

Figure 1: Prompt Template

```

# Prepare the prompts for evaluation
prompt_template = (
    f"این متن را بدون بایاس جنسیتی کن و بایاس جنسیتی آقا از بین ببر"
)

```

```

1 import datasets
2
3 # Load the test dataset
4 test_dataset = datasets.load_dataset("AmirMohammadFakhimi/gender_neutralize", split="test")

```

```

5
6 # Prepare the prompts for evaluation
7 prompt_template = (
8     *** "
9 )
10
11 def prepare_prompt(sample):
12     return {
13         "prompt": prompt_template.format(dialog=sample["dialogue"]),
14         "summary": sample["summary"]
15     }
16
17 test_dataset = test_dataset.map(prepare_prompt, remove_columns=list(test_dataset.features))

```

Explanations:

First, import the necessary library for handling datasets.

```
1 import datasets
```

Explanations:

Load the test dataset from the Hugging Face dataset repository.

```
1 test_dataset = datasets.load_dataset("AmirMohammadFakhimi/gender_neutralize", split="test")
```

Explanations:

Prepare a prompt template for evaluating the text generation. The template asks to neutralize gender bias in the provided dialogue.

```
1 prompt_template = (
2     *** "
3 )
```

Explanations:

Define a function to prepare prompts for each sample in the dataset.

```
1 def prepare_prompt(sample):
2     return {
3         "prompt": prompt_template.format(dialog=sample["dialogue"]),
4         "summary": sample["summary"]
5     }
```

Explanations:

Apply the 'prepare_prompt' function to the dataset and remove original columns.

```
1 test_dataset = test_dataset.map(prepare_prompt, remove_columns=list(test_dataset.features))
```

```

1 import torch
2
3 def generate_predictions(model, tokenizer, dataset):
4     predictions = []
5     references = []
6
7     model.eval()
8     with torch.no_grad():
9         for sample in dataset:

```

```

10 model_input = tokenizer(sample["prompt"], return_tensors="pt").to("cuda")
11 output = model.generate(**model_input, max_new_tokens=100)
12 prediction = tokenizer.decode(output[0], skip_special_tokens=True)
13
14 # Extract the part after "*** :"
15 prediction = prediction.split('\n_____\n***:\n')[1].strip()
16 prediction = re.sub('***:', '', prediction)
17 # Remove non-Alphabetical characters and digits except dot
18 prediction = re.sub(r'[^.A-Za-z\s]', '', prediction)
19 # Define a regular expression pattern that matches one or more spaces
20 pattern = re.compile(r" +")
21 # Apply the pattern to the text and replace the matches with a single space
22 prediction = pattern.sub(" ", prediction)
23
24 predictions.append(prediction)
25 references.append(sample["summary"])
26
27 return predictions, references
28
29 # Generate predictions for base model
30 base_predictions, references = generate_predictions(model, tokenizer, test_dataset)

```

Explanations:

Import the torch library for tensor computations.

```
1 import torch
```

Explanations:

Define a function to generate predictions using the model and tokenizer. The function processes each sample in the dataset, generates predictions, and cleans the generated text.

```

1 def generate_predictions(model, tokenizer, dataset):
2     predictions = []
3     references = []
4
5     model.eval()
6     with torch.no_grad():
7         for sample in dataset:
8             model_input = tokenizer(sample["prompt"], return_tensors="pt").to("cuda")
9             output = model.generate(**model_input, max_new_tokens=100)
10            prediction = tokenizer.decode(output[0], skip_special_tokens=True)
11
12            # Extract the part after "    :"
13            prediction = prediction.split('\n_____\n***:\n')[1].strip()
14            prediction = re.sub('***:', '', prediction)
15            # Remove non-Alphabetical characters and digits except dot
16            prediction = re.sub(r'[^.A-Za-z\s]', '', prediction)
17            # Define a regular expression pattern that matches one or more spaces
18            pattern = re.compile(r" +")
19            # Apply the pattern to the text and replace the matches with a single space
20            prediction = pattern.sub(" ", prediction)
21
22            predictions.append(prediction)
23            references.append(sample["summary"])
24
25     return predictions, references

```

Explanations:

Generate predictions for the base model using the defined function.

```
1 base_predictions, references = generate_predictions(model, tokenizer, test_dataset)
```

```

1 from evaluate import load
2 from hazm import word_tokenize
3
4 bertscore_metric = load("bertscore")
5 bleu_metric = load('bleu')
6
7 def compute_metrics(predictions, references):
8     # fix the references with the BLEU score HuggingFace format
9     references = [[ref] for ref in references]
10
11     # Compute BLEU score
12     bleu_score = bleu_metric.compute(predictions=predictions, references=references, tokenizer=
13         word_tokenize)["bleu"]
14
15     # Compute BERTScore
16     bertscore = bertscore_metric.compute(predictions=predictions, references=references, lang="fa",
17         model_type='microsoft/deberta-v2-xxlarge-mnli')
18     bertscore_f1 = sum(bertscore["f1"]) / len(bertscore["f1"])
19
20     return bleu_score, bertscore_f1
21
22 # Compute metrics for both models
23 base_bleu, base_bertscore = compute_metrics(base_predictions, references)
24
25 print(f"Base Model - BLEU: {base_bleu}, BERTScore F1: {base_bertscore}")

```

Explanations:

Import the 'evaluate' library for computing evaluation metrics and 'hazm' for tokenizing Persian text.

```

1 from evaluate import load
2 from hazm import word_tokenize

```

Explanations:

Load the BERTScore and BLEU metrics.

```

1 bertscore_metric = load("bertscore")
2 bleu_metric = load('bleu')

```

Explanations:

Define a function to compute evaluation metrics for the predictions. The function computes the BLEU score and the BERTScore F1 metric.

```

1 def compute_metrics(predictions, references):
2     # fix the references with the BLEU score HuggingFace format
3     references = [[ref] for ref in references]
4
5     # Compute BLEU score
6     bleu_score = bleu_metric.compute(predictions=predictions, references=references,
7         tokenizer=word_tokenize)
8
9     # Compute BERTScore
10    bertscore = bertscore_metric.compute(predictions=predictions, references=references,
11        lang="fa", model_type='microsoft/deberta-v2-xxlarge-mnli')
12    bertscore_f1 = sum(bertscore["f1"]) / len(bertscore["f1"])
13
14    return bleu_score, bertscore_f1

```

Explanations:

Compute evaluation metrics for the base model's predictions and print the results.

```
1 base_bleu, base_bertscore = compute_metrics(base_predictions, references)
2
3 print(f"Base Model - BLEU: {base_bleu}, BERTScore F1: {base_bertscore}")
```

Step 3

Part of this code is taken from The LLaMa 3 summarization example using the samsum dataset. We changed it significantly to use our dataset and preprocess our data to be usable on LLaMa 3.

It is noteworthy that we used "packing" for batching strategy instead of "padding". It is better for data with a size of more than 1000.

Explanations:

The 'dataclass' decorator is used to create a simple class for the dataset configuration, containing the dataset name and the train and test split names.

```
1 from dataclasses import dataclass
2
3 @dataclass
4 class samsum_dataset:
5     dataset: str = "samsum_dataset"
6     train_split: str = "train"
7     test_split: str = "validation"
```

Explanations:

The 'get_preprocessed_samsum' function loads the dataset, applies a prompt template, tokenizes the text, and prepares the data for model training.

First, it loads the dataset using the 'datasets.load_dataset' method and specifies the split (train or test) to be loaded.

```
1 import copy
2 import datasets
3
4 def get_preprocessed_samsum(dataset_config, tokenizer, split):
5     dataset = datasets.load_dataset("AmirMohammadFakhimi/gender_neutralize", split=split)
```

Explanations:

A prompt template is defined to instruct the model on the task, which is to remove gender bias from the provided dialogue.

```
1 prompt = (
2     f"*** "
3 )
```

Explanations:

The 'apply_prompt_template' function is applied to each sample in the dataset to create the prompt and summary pairs.

```
1 def apply_prompt_template(sample):
2     return {
3         "prompt": prompt.format(dialog=sample["dialogue"]),
4         "summary": sample["summary"],
5     }
6
```

```
dataset = dataset.map(apply\_prompt\_template, remove\_columns=list(dataset.features))
```

Explanations:

The 'tokenize_add_label' function tokenizes the prompt and summary, concatenates them, and prepares the input IDs, attention mask, and labels for training.

```
def tokenize_add_label(sample):
    prompt = tokenizer.encode(tokenizer.bos_token + sample["prompt"], add_special_tokens=False)
    summary = tokenizer.encode(sample["summary"] + tokenizer.eos_token, add_special_tokens=False)

    sample = {
        "input_ids": prompt + summary,
        "attention_mask": [1] * (len(prompt) + len(summary)),
        "labels": [-100] * len(prompt) + summary,
    }

    return sample

dataset = dataset.map(tokenize\_add\_label, remove\_columns=list(dataset.features))
```

Explanations:

The 'get_preprocessed_dataset' function is a wrapper around 'get_preprocessed_samsum' that selects the appropriate split (train or test) based on the provided configuration and returns the preprocessed dataset.

```
def get_preprocessed_dataset(
    tokenizer, dataset_config, split: str = "train"
) -> torch.utils.data.Dataset:

    def get_split():
        return (
            dataset_config.train_split
            if split == "train"
            else dataset_config.test_split
        )

    return get_preprocessed_samsum(
        dataset_config,
        tokenizer,
        get_split(),
    )
```

Explanations:

The 'ConcatDataset' and 'get_dataloader_kwargs' functions are imported to concatenate the dataset if necessary and to get the keyword arguments for the data loader.

```
from llama_recipes.data.concatenator import ConcatDataset
from llama_recipes.utils.config_utils import get_dataloader_kwargs
```

Explanations:

The preprocessed training dataset is obtained using the 'get_preprocessed_dataset' function.

```
train_dataset = get_preprocessed_dataset(tokenizer, samsum_dataset, 'train')
```

Explanations:

The keyword arguments for the data loader are obtained using the 'get_dataloader_kwargs' function.

```
1 train_dl_kwargs = get_dataloader_kwargs(train_config, train_dataset, tokenizer, "train")
2
```

Explanations:

If the batching strategy is set to "packing," the dataset is concatenated into chunks of the specified context length.

```
1 if train_config.batching_strategy == "packing":
2     train_dataset = ConcatDataset(train_dataset, chunk_size=train_config.context_length)
```

Explanations:

Finally, a PyTorch DataLoader is created for the training dataset with the specified number of workers and other configurations.

```
1 train_dataloader = torch.utils.data.DataLoader(
2     train_dataset,
3     num_workers=train_config.num_workers_dataloader,
4     pin_memory=True,
5     **train_dl_kwargs,
6 )
```

Step 4

Preparing mode for PEFT.

```
1 from peft import get_peft_model, prepare_model_for_kbit_training, LoraConfig
2 from dataclasses import asdict
3 from llama_recipes.configs import lora_config as LORA_CONFIG
4
5 lora_config = LORA_CONFIG()
6 lora_config.r = 8
7 lora_config.lora_alpha = 32
8 lora_dropout: float=0.01
9
10 peft_config = LoraConfig(**asdict(lora_config))
11
12 model = prepare_model_for_kbit_training(model)
13 model = get_peft_model(model, peft_config)
```

Explanations:

The 'peft' library is imported, which includes functions for getting a PEFT (Parameter Efficient Fine-Tuning) model, preparing a model for k-bit training, and configuring LoRA (Low-Rank Adaptation) settings.

```
1 from peft import get_peft_model, prepare_model_for_kbit_training, LoraConfig
2 from dataclasses import asdict
3 from llama_recipes.configs import lora_config as LORA_CONFIG
```

Explanations:

The LoRA configuration is loaded and modified. Here, 'r' is the rank of the decomposition, 'lora_alpha' is the scaling factor, and 'lora_dropout' is the dropout rate.

```
1 lora_config = LORA_CONFIG()
2 lora_config.r = 8
3 lora_config.lora_alpha = 32
```

```
4 lora_dropout: float=0.01
```

Explanations:

The LoRA configuration is converted into a dictionary and used to create a 'LoraConfig' object.

```
1 peft_config = LoraConfig(**asdict(lora_config))
```

Explanations:

The model is prepared for k-bit training using the 'prepare_model_for_kbit_training' function, which optimizes the model for efficient fine-tuning. Then, the 'get_peft_model' function applies the LoRA configuration to the model.

```
1 model = prepare_model_for_kbit_training(model)
2 model = get_peft_model(model, peft_config)
```

Step 5

This is the fine tuning step.

```
1 import torch.optim as optim
2 from llama_recipes.utils.train_utils import train
3 from torch.optim.lr_scheduler import StepLR
4
5 model.train()
6
7 optimizer = optim.AdamW(
8     model.parameters(),
9     lr=train_config.lr,
10    weight_decay=train_config.weight_decay,
11 )
12 scheduler = StepLR(optimizer, step_size=1, gamma=train_config.gamma)
13
14 # Start the training process
15 results = train(
16     model,
17     train_dataloader,
18     None,
19     tokenizer,
20     optimizer,
21     scheduler,
22     train_config.gradient_accumulation_steps,
23     train_config,
24     None,
25     None,
26     None,
27     wandb_run=None,
28 )
```

Explanations:

The 'torch.optim' module is imported to handle the optimization process. The 'train' function and 'StepLR' scheduler are imported to manage training and learning rate scheduling, respectively.

```
1 import torch.optim as optim
2 from llama_recipes.utils.train_utils import train
3 from torch.optim.lr_scheduler import StepLR
```


Explanations:

The model is set to training mode using 'model.train()'. This ensures that layers like dropout are applied appropriately during training.

```
1 model.train()
```

Explanations:

An AdamW optimizer is created with parameters from the model, learning rate from the training configuration, and weight decay to regularize the model.

```
1 optimizer = optim.AdamW(  
2     model.parameters(),  
3     lr=train_config.lr,  
4     weight_decay=train_config.weight_decay,  
5 )
```

Explanations:

A learning rate scheduler, 'StepLR', is set up to adjust the learning rate after every epoch. The 'step_size' parameter defines the number of epochs after which the learning rate is updated, and 'gamma' is the multiplicative factor for learning rate decay.

```
1 scheduler = StepLR(optimizer, step_size=1, gamma=train_config.gamma)
```

Explanations:

The 'train' function is called to start the training process. It takes the model, data loader, tokenizer, optimizer, scheduler, gradient accumulation steps, and other training configurations as inputs. 'wandb_run' is set to 'None' as no Weights & Biases run is being used here.

```
1 # Start the training process  
2 results = train(  
3     model,  
4     train_dataloader,  
5     None,  
6     tokenizer,  
7     optimizer,  
8     scheduler,  
9     train_config.gradient_accumulation_steps,  
10    train_config,  
11    None,  
12    None,  
13    None,  
14    wandb_run=None,  
15 )
```

Step 6

```
1 model.save_pretrained(train_config.output_dir)
```

Explanations:

The trained model is saved to the directory specified in the training configuration's 'output_dir'. This allows the model to be reloaded and used later without retraining.

```
1 model.save_pretrained(train_config.output_dir)
```

```

1 from transformers import LlamaForCausalLM, AutoTokenizer, BitsAndBytesConfig
2 from peft import PeftModel, PeftConfig
3 import torch
4
5 # Load the base model and tokenizer
6 base_model_name = "meta-llama/Meta-Llama-3-8B"
7 tokenizer = AutoTokenizer.from_pretrained(base_model_name)
8 tokenizer.pad_token = tokenizer.eos_token
9
10 # Use mixed precision and quantization to save memory
11 config = BitsAndBytesConfig(
12     load_in_8bit=True,
13     llm_int8_enable_fp32_cpu_offload=True # Enable CPU offload to save GPU memory
14 )
15
16 base_model = LlamaForCausalLM.from_pretrained(
17     base_model_name,
18     device_map="auto",
19     quantization_config=config,
20     torch_dtype=torch.float16
21 )
22
23 # Load the adapter configuration
24 peft_model_id = '/kaggle/input/meta-llama-sexbias/pytorch/default/1'
25 adapter_config = PeftConfig.from_pretrained(peft_model_id)
26
27 # Load the LoRA adapter
28 model = PeftModel.from_pretrained(base_model, peft_model_id, config=adapter_config)

```

Explanations:

The necessary libraries for loading and configuring the model for inference are imported. These include 'LlamaForCausalLM', 'AutoTokenizer', and 'BitsAndBytesConfig' from 'transformers', as well as 'PeftModel' and 'PeftConfig' from 'peft'.

```

1 from transformers import LlamaForCausalLM, AutoTokenizer, BitsAndBytesConfig
2 from peft import PeftModel, PeftConfig
3 import torch

```

Explanations:

The base model name is defined, and the tokenizer is loaded using the 'AutoTokenizer' class. The end-of-sequence token is set as the padding token.

```

1 base_model_name = "meta-llama/Meta-Llama-3-8B"
2 tokenizer = AutoTokenizer.from_pretrained(base_model_name)
3 tokenizer.pad_token = tokenizer.eos_token

```

Explanations:

A 'BitsAndBytesConfig' is created to use mixed precision and quantization, which helps save memory by loading the model in 8-bit precision and enabling FP32 CPU offloading.

```

1 config = BitsAndBytesConfig(
2     load_in_8bit=True,
3     llm_int8_enable_fp32_cpu_offload=True # Enable CPU offload to save GPU memory
4 )

```

Explanations:

The base model is loaded using the 'LlamaForCausalLM' class with the defined quantization configuration and mixed precision (FP16).

```

1 base_model = LlamaForCausalLM.from_pretrained(
2     base_model_name,
3     device_map="auto",
4     quantization_config=config,
5     torch_dtype=torch.float16
6 )

```

Explanations:

The adapter configuration is loaded from a specified directory using the 'PeftConfig' class. This configuration contains settings specific to the PEFT model.

```

1 peft_model_id = '/kaggle/input/meta-llama-sexbias/pytorch/default/1'
2 adapter_config = PeftConfig.from_pretrained(peft_model_id)

```

Explanations:

The LoRA adapter is loaded using the 'PeftModel' class, which combines the base model with the PEFT configuration.

```

1 model = PeftModel.from_pretrained(base_model, peft_model_id, config=adapter_config)

```

Step 7

```

1 import re
2
3 def generate_predictions(model, tokenizer, dataset):
4     predictions = []
5     references = []
6
7     model.eval()
8     with torch.no_grad():
9         for sample in dataset:
10             model_input = tokenizer(sample["prompt"], return_tensors="pt").to("cuda")
11             output = model.generate(**model_input, max_new_tokens=100)
12             prediction = tokenizer.decode(output[0], skip_special_tokens=True)
13
14             # Extract the part after "***"
15             prediction = prediction.split('***')[1].strip()
16             prediction = re.sub('***:', ' ', prediction)
17             # Remove non-Alphabetical characters and digits except dot
18             prediction = re.sub(r'[^-A-Za-z\s]', ' ', prediction)
19             # Define a regular expression pattern that matches one or more spaces
20             pattern = re.compile(r" +")
21             # Apply the pattern to the text and replace the matches with a single space
22             prediction = pattern.sub(" ", prediction)
23
24             predictions.append(prediction)
25             references.append(sample["summary"])
26
27             # Clear CUDA cache to free up memory
28             torch.cuda.empty_cache()
29
30     return predictions, references

```

Explanations:

The 'generate_predictions' function generates predictions from the model for each sample in the dataset and cleans up the text.

- The function sets the model to evaluation mode using 'model.eval()' and disables gradient calculations

with 'torch.no_grad()':

- For each sample in the dataset, it tokenizes the input prompt and generates a prediction using the model.
- The predicted text is decoded and cleaned:
 - It extracts the relevant portion of the prediction after the marker "Text without gender bias" (in Persian).
 - Non-alphabetical characters and digits are removed, except for dots.
 - Multiple spaces are replaced with a single space.
- Predictions and references are collected, and CUDA memory is cleared with 'torch.cuda.empty_cache()' to manage GPU resources.

```
1 def generate_predictions(model, tokenizer, dataset):
2     predictions = []
3     references = []
4
5     model.eval()
6     with torch.no_grad():
7         for sample in dataset:
8             model_input = tokenizer(sample["prompt"], return_tensors="pt").to("cuda")
9             output = model.generate(**model_input, max_new_tokens=100)
10            prediction = tokenizer.decode(output[0], skip_special_tokens=True)
11
12            # Extract the part after "***"
13            prediction = prediction.split('***')[1].strip()
14            prediction = re.sub('***', '', prediction)
15            # Remove non-Alphabetical characters and digits except dot
16            prediction = re.sub(r'[^.-A-Za-z\s]', '', prediction)
17            # Define a regular expression pattern that matches one or more spaces
18            pattern = re.compile(r" +")
19            # Apply the pattern to the text and replace the matches with a single space
20            prediction = pattern.sub(" ", prediction)
21
22            predictions.append(prediction)
23            references.append(sample["summary"])
24
25            # Clear CUDA cache to free up memory
26            torch.cuda.empty_cache()
27
28    return predictions, references
```

Explanations:

Predictions and references are saved to text files for later evaluation. Each line of the prediction file corresponds to a generated output, while each line of the references file corresponds to the reference summary.

```
1 with open("tuned_predictions.txt", "w") as f:
2     for pred in tuned_predictions:
3         f.write(str(pred) + "\n")
4
5 with open("references.txt", "w") as f:
6     for ref in references:
7         f.write(str(ref) + "\n")
```

Explanations:

The saved text files are read back into lists. Each line is converted to an integer. This step may be a mistake as predictions and references are typically strings. Ensure this conversion is appropriate for your use case.

```

1 references = []
2 with open("references.txt", "r") as f:
3     for line in f:
4         references.append(int(line.strip()))
5
6 tuned_predictions = []
7 with open("tuned_predictions.txt", "r") as f:
8     for line in f:
9         tuned_predictions.append(int(line.strip()))

```

Explanations:

The 'compute_metrics' function computes BLEU and BERTScore metrics for the predictions against the references. BLEU measures the precision of n-grams, while BERTScore evaluates the semantic similarity using BERT embeddings.

- BLEU score is calculated using the 'bleu_metric' with tokenized references and predictions.
- BERTScore is computed using the 'bertscore_metric' with the specified language and model.

```

1 from evaluate import load
2 from hazm import word_tokenize
3
4 bertscore_metric = load("bertscore")
5 bleu_metric = load('bleu')
6
7 def compute_metrics(predictions, references):
8     # fix the references with the BLEU score HuggingFace format
9     references = [[ref] for ref in references]
10
11     # Compute BLEU score
12     bleu_score = bleu_metric.compute(predictions=predictions, references=references,
13                                     tokenizer=word_tokenize)
14
15     # Compute BERTScore
16     bertscore = bertscore_metric.compute(predictions=predictions, references=references,
17                                         lang="fa", model_type='microsoft/deberta-v2-xxlarge-mnli')
18     bertscore_f1 = sum(bertscore["f1"]) / len(bertscore["f1"])
19
20     return bleu_score, bertscore_f1

```

Explanations:

The computed BLEU and BERTScore metrics are printed. These metrics provide an evaluation of the models performance in terms of both n-gram precision and semantic similarity.

```

1 # Compute metrics for fine-tuned model
2 tuned_bleu, tuned_bertscore = compute_metrics(tuned_predictions, references)
3
4 print(f"Fine-tuned Model - BLEU: {tuned_bleu}, BERTScore F1: {tuned_bertscore}")

```

Results

In the following section, we have included some examples of the results and the overall metrics of the model. An important note is that for some of the results, the not-fine-tunes model wanted to repeat the answer multiple times. In those cases, we added post-processing to remove the duplications of text and used it to calculate metrics.

As shown in the figures below, we saw a huge improvement in BLUE and BERT scores after fine-tuning.

Base Model - BLEU: {'bleu': 0.19212972085620741, 'precisions': [0.2688223938223938, 0.2023511755877939, 0.17159916926272067, 0.145979492714517], 'brevity_penalty': 1.0, 'length_ratio': 2.5009052504526252, 'translation_length': 4144, 'reference_length': 1657}, BERTScore F1: 0.818024289730477

Figure 2: Base model Evaluation Metrics

Fine-tuned Model - BLEU: {'bleu': 0.4282493906811996, 'precisions': [0.6205770277626566, 0.46599645180366644, 0.37928802588996763, 0.30664760543245173], 'brevity_penalty': 1.0, 'length_ratio': 1.1086300543150271, 'translation_length': 1837, 'reference_length': 1657}, BERTScore F1: 0.9234269032739613

Figure 3: Fine-tuned model Evaluation Metrics

این متن را بدون بایاس جنسیتی کن و بایاس جنسیتی آنرا از بین ببر
مردان باید در برنامه‌ریزی برای تعطیلات خانوادگی پیش‌قدم باشند

متن بدون بایاس جنسیتی:
شما باید در برنامه‌ریزی برای تعطیلات خانوادگی پیش‌قدم باشید

(a) Base Model

این متن را بدون بایاس جنسیتی کن و بایاس جنسیتی آنرا از بین ببر
مردان باید در برنامه‌ریزی برای تعطیلات خانوادگی پیش‌قدم باشند

متن بدون بایاس جنسیتی:
برنامه‌ریزی برای تعطیلات خانوادگی باید به عهده‌ی فردی باشد که توانایی‌های لازم را دارد

(b) Fine-Tuned Model

Figure 4: Sample Prediction 1

:این متن را بدون بایاس جنسیتی کن و بایاس جنسیتی آنرا از بین ببر
افراد ممکن است به جنبه های مختلف زندگی خود توجه کنند و تمرکز متفاوتی بر روی مهارت های حرفه ای داشته باشند

:متن بدون بایاس جنسیتی

افراد ممکن است به جنبه های مختلف زندگی خود توجه کنند و تمرکز متفاوتی بر روی مهارت های حرفه ای داشته باشند

:متن بدون بایاس جنسیتی

افراد ممکن است به جنبه های مختلف زندگی خود توجه کنند و تمرکز متفاوتی بر روی مهارت های حرفه ای داشته باشند

:متن بدون بایاس جنسیتی

افراد ممکن است به جنبه های مختلف زندگی خود توجه کنند و تمرکز متفاوتی بر روی مه

(a) Base Model

:این متن را بدون بایاس جنسیتی کن و بایاس جنسیتی آنرا از بین ببر
افراد ممکن است به جنبه های مختلف زندگی خود توجه کنند و تمرکز متفاوتی بر روی مهارت های حرفه ای داشته باشند

:متن بدون بایاس جنسیتی

افراد ممکن است به جنبه های مختلف زندگی خود توجه کنند و تمرکز متفاوتی بر روی مهارت های حرفه ای داشته باشند

(b) Fine-Tuned Model

Figure 5: Sample Prediction 2

:این متن را بدون بایاس جنسیتی کن و بایاس جنسیتی آنرا از بین ببر
زنان باید سرپرست سازماندهی رویدادهای اجتماعی باشند

:متن بدون بایاس جنسیتی

زنان و مردان باید سرپرست سازماندهی رویدادهای اجتماعی باشند

(a) Base Model

:این متن را بدون بایاس جنسیتی کن و بایاس جنسیتی آنرا از بین ببر
زنان باید سرپرست سازماندهی رویدادهای اجتماعی باشند

:متن بدون بایاس جنسیتی

سرپرستی سازماندهی رویدادهای اجتماعی می تواند بر اساس توانایی های فردی و علاقه های شخصی تقسیم شود

(b) Fine-Tuned Model

Figure 6: Sample Prediction 3

:این متن را بدون بایاس جنسیتی کن و بایاس جنسیتی آنرا از بین ببر
المیرا اسم یک دختر است

:متن بدون بایاس جنسیتی
المیرا اسم یک دختر است

:متن با بایاس جنسیتی
المیرا اسم یک دختر است

:متن بدون بایاس جنسیتی
المیرا اسم یک دختر است

:متن با بایاس جنسیتی
المیرا اسم یک دختر است

:متن بدون بایاس جنسیتی
المیرا اسم یک دختر است

:متن با بایاس جنسیتی
المیرا اسم یک دختر است

(a) Base Model

:این متن را بدون بایاس جنسیتی کن و بایاس جنسیتی آنرا از بین ببر
المیرا اسم یک دختر است

:متن بدون بایاس جنسیتی
المیرا اسم یک دختر است

(b) Fine-Tuned Model

Figure 7: Sample Prediction 4

:این متن را بدون بایاس جنسیتی کن و بایاس جنسیتی آنرا از بین ببر
زن‌ها نمی‌توانند به خوبی مردها طراحی کنند.

:متن بدون بایاس جنسیتی

زن‌ها می‌توانند به خوبی مردها طراحی کنند.

:متن با بایاس جنسیتی

زن‌ها نمی‌توانند به خوبی مردها طراحی کنند.

:متن با بایاس جنسیتی

زن‌ها نمی‌توانند به خوبی مردها طراحی کنند.

:متن با بایاس جنسیتی

زن‌ها نمی‌توانند به خوبی مردها طراحی کنند.

:متن با بایاس جنسیتی

زن‌ها نمی‌توانند به خوبی مردها طراحی کنند.

متن با

(a) Base Model

:این متن را بدون بایاس جنسیتی کن و بایاس جنسیتی آنرا از بین ببر
زن‌ها نمی‌توانند به خوبی مردها طراحی کنند.

:متن بدون بایاس جنسیتی

هرکسی می‌تواند به خوبی طراحی کند.

(b) Fine-Tuned Model

Figure 8: Sample Prediction 5

:این متن را بدون بایاس جنسیتی کن و بایاس جنسیتی آنرا از بین ببر
زن‌ها در ورزش‌های سنگین موفق نمی‌شوند

:متن بدون بایاس جنسیتی
زن‌ها در ورزش‌های سنگین موفق نمی‌شوند

(a) Base Model

:این متن را بدون بایاس جنسیتی کن و بایاس جنسیتی آنرا از بین ببر
زن‌ها در ورزش‌های سنگین موفق نمی‌شوند

:متن بدون بایاس جنسیتی
موفقیت در ورزش‌های سنگین به توانایی‌های فیزیکی و ذهنی فرد بستگی دارد

(b) Fine-Tuned Model

Figure 9: Sample Prediction 6