

```
1 #%% md
2 # Deep Learning Assignment 2 Part 1: Image
3 # captioning with RNNs and LSTMs
4 #%% md
5 # Image Captioning with RNNs
6 In this exercise, you will implement vanilla
recurrent neural networks (RNNs), [long-short term
memory networks (LSTMs)](https://www.researchgate.net/publication/13853244\_Long\_Short-term\_Memory),
and [attention-based LSTMs](https://arxiv.org/abs/1409.0473) to train a model that can generate
natural language captions for images.
7
8 Models in this exercise are highly similar to very
early works in neural-network based image
captioning. If you are interested to learn more,
check out these two papers:
9
10 1. [Show and Tell: A Neural Image Caption Generator
11 ](https://arxiv.org/abs/1411.4555)
12 2. [Show, Attend and Tell: Neural Image Caption
13 Generation with Visual Attention](https://arxiv.org/abs/1502.03044)
14 #%% md
15 ## Setup Code
16
17 Before getting started, we need to run some
boilerplate code to set up our environment, same as
previous assignments. You'll need to rerun this
setup code each time you start the notebook.
18
19 First, run this cell load the [autoreload](https://ipython.readthedocs.io/en/stable/config/extensions/autoreload.html?highlight=autoreload) extension.
This allows us to edit .py source files, and re-
import them into the notebook for a seamless
editing and debugging experience.
20 #%%
```

```
21 #%% md
22 ### Google Colab Setup
23
24 Next we need to run a few commands to set up our
   environment on Google Colab. If you are running
   this notebook on a local machine you can skip this
   section.
25
26 Run the following cell to mount your Google Drive.
   Follow the link, sign in to your Google account (
   the same account you used to store this notebook!)
   and copy the authorization code into the text box
   that appears below.
27 #%%
28 # from google.colab import drive
29 #
30 # drive.mount("/content/drive")
31 #%% md
32 Now recall the path in your Google Drive where you
   uploaded this notebook, fill it in below. If
   everything is working correctly then running the
   following cell should print the filenames from the
   assignment:
33
34 ``
35 ["dl", "a2_helper.py", "rnn_lstm_captioning.ipynb",
   "rnn_lstm_captioning.py", "transformers.py", "Transformers.ipynb", "two_digit_op.json"]
36 ``
37 #%%
38 # import os
39 # import sys
40 #
41 # # TODO: Enter the foldername in your Drive where
   you have saved the unzipped
42 # # assignment folder, e.g. 'dl/assignments/
   assignment2/'
43 # FOLDERNAME = 'dl/assignments/assignment2/'
44 # assert FOLDERNAME is not None, "[!] Enter the
   foldername."
45 #
```

```
46 # # Now that we've mounted your Drive, this ensures
    # that
47 # # the Python interpreter of the Colab VM can load
48 # # python files from within it.
49 # import sys
50 # sys.path.append('/content/drive/My Drive/{}'.format(FOLDERNAME))
51 #%% md
52 Once you have successfully mounted your Google
Drive and located the path to this assignment, run
the following cell to allow us to import from the
`.py` files of this assignment. If it works
correctly, it should print the message:
53
54 ``
55 Hello from rnn_lstm_captioning.py!
56 ``
57
58 as well as the last edit time for the file `rnn_lstm_captioning.py`.
59 #%%
60 import os
61 import time
62 from pathlib import Path
63 from rnn_lstm_captioning import
    hello_rnn_lstm_captioning
64
65
66 os.environ["TZ"] = "Asia/Jerusalem"
67 time.tzset()
68 NOTEBOOK_DIR = Path.cwd()
69 candidates = [NOTEBOOK_DIR] + list(NOTEBOOK_DIR.
    parents)
70 ASSIGNMENT_ROOT = None
71 for p in candidates:
72     if (p / "rnn_lstm_captioning.py").exists() and
        (p / "a2_helper.py").exists():
73         ASSIGNMENT_ROOT = p
74         break
75 print("Assignment root:", ASSIGNMENT_ROOT)
76 print("Files:", sorted([p.name for p in
```

```
76 ASSIGNMENT_ROOT.iteaddir()])[:20])
77
78 hello_rnn_lstm_captioning()
79
80 rnn_lstm_path = Path(__file__).parent / "
81     rnn_lstm_captioning.py" if "__file__" in globals
82     () else (Path.cwd() / "rnn_lstm_captioning.py")
83 if not rnn_lstm_path.exists():
84     # fallback: use the ASSIGNMENT_ROOT we defined
85     # earlier
86     rnn_lstm_path = ASSIGNMENT_ROOT / "
87         rnn_lstm_captioning.py"
88
89 print("Using:", rnn_lstm_path)
90 print("rnn_lstm_captioning.py last edited on",
91       time.ctime(rnn_lstm_path.stat().st_mtime))
92 #%% md
93 ### Load Packages
94
95 Run some setup code for this notebook: Import some
96     useful packages and increase the default figure
97     size.
98
99 #%%
100 import math
101 import os
102 import time
103 import matplotlib.pyplot as plt
104 import torch
105 from torch import nn
106
107 from dl.grad import compute_numeric_gradient,
108     rel_error
109 from dl.utils import attention_visualizer,
110     reset_seed
111
112 # for plotting
113 %matplotlib inline
114 plt.style.use("seaborn-v0_8") # Prettier plots
115 plt.rcParams["figure.figsize"] = (10.0, 8.0) #  
set default size of plots
```

```
107 plt.rcParams["font.size"] = 24
108 plt.rcParams["image.interpolation"] = "nearest"
109 plt.rcParams["image.cmap"] = "gray"
110 %% md
111 We will use GPUs to accelerate our computation in
this notebook. Run the following to make sure GPUs
are enabled:
112 %%
113 # Device selection logic (unchanged)
114 if torch.backends.mps.is_available() and torch.
    backends.mps.is_built():
115     DEVICE = torch.device("mps")
116     print("Using Apple Silicon GPU (MPS)")
117 elif torch.cuda.is_available():
118     DEVICE = torch.device("cuda")
119     print("Using CUDA GPU")
120 else:
121     DEVICE = torch.device("cpu")
122     print("Using CPU")
123
124
125 # Define some common variables for dtypes/devices.
126 # These can be keyword arguments while defining
new tensors.
127 # NOTE: Apple MPS does NOT support float64. We
keep float32 on GPU,
128 # but force float64 tensors (used in numeric
gradient checks) onto CPU.
129
130 to_float = {"dtype": torch.float32, "device":
DEVICE}
131
132 DOUBLE_DEVICE = DEVICE
133 if DEVICE.type == "mps":
134     DOUBLE_DEVICE = torch.device("cpu")
135
136 to_double = {"dtype": torch.float64, "device":
DOUBLE_DEVICE}
137 %% md
138 # COCO Captions
139
```

- 140 For this exercise we will use the 2014 release of the [COCO Captions dataset](<http://cocodataset.org/>) which has become the standard testbed for image captioning. The dataset consists of 80,000 training images and 40,000 validation images, each annotated with 5 captions written by workers on Amazon Mechanical Turk.
- 141
- 142 We have preprocessed the data for you already and saved them into a serialized data file. It contains 10,000 image-caption pairs for training and 500 for testing. The images have been downsampled to 112x112 for computation efficiency and captions are tokenized and numericalized, clamped to 15 words. You can download the file named `coco.pt` (378MB) with the link below and run some useful stats.
- 143
- 144 You will later use RegNet-X 400MF model to extract features for the images. A few notes on the caption preprocessing:
- 145
- 146 Dealing with strings is inefficient, so we will work with an encoded version of the captions. Each word is assigned an integer ID, allowing us to represent a caption by a sequence of integers. The mapping between integer IDs and words is saved in an entry named `vocab` (both `idx_to_token` and `token_to_idx`), and we use the function `decode_captions` from `a2_helper.py` to convert tensors of integer IDs back into strings.
- 147
- 148 There are a couple special tokens that we add to the vocabulary. We prepend a special `` token and append an `` token to the beginning and end of each caption respectively. Rare words are replaced with a special `` token (for "unknown"). In addition, since we want to train with minibatches containing captions of different lengths, we pad short captions with a special `<NULL>` token after the `` token and don't

```
148 compute loss or gradient for `<NULL>` tokens.  
    Since they are a bit of a pain, we have taken care  
    of all implementation details around special  
    tokens for you.  
149 #%%  
150 import multiprocessing  
151  
152 # Set a few constants related to data loading.  
153 IMAGE_SHAPE = (112, 112)  
154 NUM_WORKERS = multiprocessing.cpu_count()  
155  
156 # Batch size used for full training runs:  
157 BATCH_SIZE = 256  
158  
159 # Batch size used for overfitting sanity checks:  
160 OVR_BATCH_SIZE = BATCH_SIZE // 8  
161  
162 # Batch size used for visualization:  
163 VIS_BATCH_SIZE = 4  
164 #%%  
165 import urllib.request  
166 from a2_helper import load_coco_captions  
167  
168 # Download and load serialized COCO data from coco  
# .pt  
169 # It contains a dictionary of  
170 # "train_images" - resized training images (  
# IMAGE_SHAPE)  
171 # "val_images" - resized validation images (  
# IMAGE_SHAPE)  
172 # "train_captions" - tokenized and numericalized  
# training captions  
173 # "val_captions" - tokenized and numericalized  
# validation captions  
174 # "vocab" - caption vocabulary, including "  
# idx_to_token" and "token_to_idx"  
175  
176 os.makedirs("./datasets", exist_ok=True)  
177 dst = "./datasets/coco.pt"  
178  
179 if os.path.isfile(dst):
```

```
180     print("COCO data exists!")
181 else:
182     print("Downloading COCO dataset...")
183     url = "http://web.eecs.umich.edu/~justincj/
184         teaching/eecs498/coco.pt"
185     urllib.request.urlretrieve(url, dst)
186     print("Downloaded to", dst)
187 # load COCO data from coco.pt, load_COCO is
188 # implemented in a2_helper.py
189 data_dict = load_coco_captions(path='./datasets/
190     coco.pt')
191
192 num_train = data_dict["train_images"].size(0)
193 num_val = data_dict["val_images"].size(0)
194
195 # declare variables for special tokens
196 NULL_index = data_dict["vocab"]["token_to_idx"]["<
197     NULL>"]
198 START_index = data_dict["vocab"]["token_to_idx"]["<
199     <START>>"]
200 END_index = data_dict["vocab"]["token_to_idx"]["<
201     <END>>"]
202 UNK_index = data_dict["vocab"]["token_to_idx"]["<
203     UNK>>"]
204
205 #%% md
206 ## Look at the data
207 It is always a good idea to look at examples from
208 the dataset before working with it.
209
210 Run the following to sample a small minibatch of
211 training data and show the images and their
212 captions. Running it multiple times and looking at
213 the results helps you to get a sense of the
214 dataset.
215
216 Note that we decode the captions using the `decode_
217 captions` function.
218 You can check its implementation in `a2_helper.py
219     `!
220
221 #%%
```

```
207 from a2_helper import decode_captions
208
209
210 # Sample a minibatch and show the reshaped 112x112
   images and captions
211 sample_idx = torch.randint(0, num_train, (
    VIS_BATCH_SIZE, ))
212 sample_images = data_dict["train_images"][
    sample_idx]
213 sample_captions = data_dict["train_captions"][
    sample_idx]
214 for i in range(VIS_BATCH_SIZE):
215     plt.imshow(sample_images[i].permute(1, 2, 0))
216     plt.axis("off")
217     caption_str = decode_captions(
218         sample_captions[i], data_dict["vocab"]["
    idx_to_token"])
219     )
220     plt.title(caption_str)
221     plt.show()
222 %% md
223 # Recurrent Neural Networks
224 As discussed in lecture, we will use Recurrent
   Neural Network (RNN) language models for image
   captioning. We will cover the vanilla RNN model
   first and later LSTM and attention-based language
   models.
225 %% md
226 ## Vanilla RNN: step forward
227
228 First implement the `rnn_step_forward` for a
   single timestep of a vanilla recurrent neural
   network.
229 Run the following to check your implementation.
   You should see errors on the order of `1e-8` or
   less.
230 %%
231 from rnn_lstm_captioning import rnn_step_forward
232
233 N, D, H = 3, 10, 4
234
```

```

235 x = torch.linspace(-0.4, 0.7, steps=N * D, **  

236 to_double).view(N, D)  

237 prev_h = torch.linspace(-0.2, 0.5, steps=N * H, **  

238 to_double).view(N, H)  

239 Wx = torch.linspace(-0.1, 0.9, steps=D * H, **  

240 to_double).view(D, H)  

241 Wh = torch.linspace(-0.3, 0.7, steps=H * H, **  

242 to_double).view(H, H)  

243 b = torch.linspace(-0.2, 0.4, steps=H, **to_double  

244 )  

245  

246 print(x.shape)  

247 print(prev_h.shape)  

248 print(Wx.shape)  

249 print(Wh.shape)  

250 print(b.shape)  

251 next_h, _ = rnn_step_forward(x, prev_h, Wx, Wh, b)  

252 expected_next_h = torch.tensor(  

253     [  

254         [-0.58172089, -0.50182032, -0.41232771, -0.  

255         .31410098],  

256         [0.66854692, 0.79562378, 0.87755553, 0.  

257         92795967],  

258         [0.97934501, 0.99144213, 0.99646691, 0.  

259         99854353],  

260     ],  

261     **to_double  

262 )
263  

264 print("next_h error: ", rel_error(expected_next_h,  

265 , next_h))  

266 #%% md  

267 ## Vanilla RNN: step backward  

268 Then implement the `rnn_step_backward` for a  

269 single timestep of a vanilla recurrent neural  

270 network. Run the following to numerically gradient  

271 check your implementation. You should see errors  

272 on the order of `1e-8` or less.  

273 #%%  

274 from rnn_lstm_captioning import rnn_step_backward  

275

```

```
263
264 reset_seed(0)
265
266 N, D, H = 4, 5, 6
267 x = torch.randn(N, D, **to_double)
268 h = torch.randn(N, H, **to_double)
269 Wx = torch.randn(D, H, **to_double)
270 Wh = torch.randn(H, H, **to_double)
271 b = torch.randn(H, **to_double)
272
273 out, cache = rnn_step_forward(x, h, Wx, Wh, b)
274
275 dnext_h = torch.randn(*out.shape, **to_double)
276
277 fx = lambda x: rnn_step_forward(x, h, Wx, Wh, b)[0]
278 fh = lambda h: rnn_step_forward(x, h, Wx, Wh, b)[0]
279 fWx = lambda Wx: rnn_step_forward(x, h, Wx, Wh, b)[0]
280 fWh = lambda Wh: rnn_step_forward(x, h, Wx, Wh, b)[0]
281 fb = lambda b: rnn_step_forward(x, h, Wx, Wh, b)[0]
282
283 dx_num = compute_numeric_gradient(fx, x, dnext_h)
284 dprev_h_num = compute_numeric_gradient(fh, h,
dnext_h)
285 dWx_num = compute_numeric_gradient(fWx, Wx,
dnext_h)
286 dWh_num = compute_numeric_gradient(fWh, Wh,
dnext_h)
287 db_num = compute_numeric_gradient(fb, b, dnext_h)
288
289 # YOUR_TURN: Implement rnn_step_backward
290 dx, dprev_h, dWx, dWh, db = rnn_step_backward(
dnext_h, cache)
291
292 print("dx error: ", rel_error(dx_num, dx))
293 print("dprev_h error: ", rel_error(dprev_h_num,
dprev_h))
```

```

294 print("dWx error: ", rel_error(dWx_num, dWx))
295 print("dWh error: ", rel_error(dWh_num, dWh))
296 print("db error: ", rel_error(db_num, db))
297 #%% md
298 ## Vanilla RNN: forward
299 Now that you have implemented the forward and
    backward passes for a single timestep of a vanilla
    RNN, you will combine these pieces to implement a
    RNN that processes an entire sequence of data.
    First implement `rnn_forward` by making calls to
    the `rnn_step_forward` function that you defined
    earlier.
300
301 Run the following to check your implementation.
    You should see errors on the order of `1e-6` or
    less.
302
303
304 #%%
305 from rnn_lstm_captioning import rnn_forward
306
307
308 N, T, D, H = 2, 3, 4, 5
309
310 x = torch.linspace(-0.1, 0.3, steps=N * T * D, **  

    to_double).view(N, T, D)
311 h0 = torch.linspace(-0.3, 0.1, steps=N * H, **  

    to_double).view(N, H)
312 Wx = torch.linspace(-0.2, 0.4, steps=D * H, **  

    to_double).view(D, H)
313 Wh = torch.linspace(-0.4, 0.1, steps=H * H, **  

    to_double).view(H, H)
314 b = torch.linspace(-0.7, 0.1, steps=H, **to_double  

    )
315
316 # YOUR_TURN: Implement rnn_forward
317 h, _ = rnn_forward(x, h0, Wx, Wh, b)
318 expected_h = torch.tensor(  

319     [
320         [  

321             [-0.42070749, -0.27279261, -0.11074945

```

```

321 , 0.05740409, 0.22236251],
322 [-0.39525808, -0.22554661, -0.0409454
, 0.14649412, 0.32397316],
323 [-0.42305111, -0.24223728, -0.04287027
, 0.15997045, 0.35014525],
324 ],
325 [
326 [-0.55857474, -0.39065825, -0.19198182
, 0.02378408, 0.23735671],
327 [-0.27150199, -0.07088804, 0.13562939
, 0.33099728, 0.50158768],
328 [-0.51014825, -0.30524429, -0.06755202
, 0.17806392, 0.40333043],
329 ],
330 ],
331 **to_double
332 )
333 print("h error: ", rel_error(expected_h, h))
334 #%% md
335 ## Vanilla RNN: backward
336 Implement the `rnn_backward` for a vanilla RNN.
This should run back-propagation over the entire
sequence, making calls to the `rnn_step_backward`
function that you defined earlier. You should see
errors on the order of `1e-6` or less.
337
338 #%%
339 from rnn_lstm_captioning import rnn_backward,
rnn_forward
340
341 reset_seed(0)
342
343 N, D, T, H = 2, 3, 10, 5
344
345 x = torch.randn(N, T, D, **to_double)
346 h0 = torch.randn(N, H, **to_double)
347 Wx = torch.randn(D, H, **to_double)
348 Wh = torch.randn(H, H, **to_double)
349 b = torch.randn(H, **to_double)
350
351 out, cache = rnn_forward(x, h0, Wx, Wh, b)

```

```

352
353 dout = torch.randn(*out.shape, **to_double)
354
355 # YOUR_TURN: Implement rnn_backward
356 dx, dh0, dWx, dWh, db = rnn_backward(dout, cache)
357
358 fx = lambda x: rnn_forward(x, h0, Wx, Wh, b)[0]
359 fh0 = lambda h0: rnn_forward(x, h0, Wx, Wh, b)[0]
360 fWx = lambda Wx: rnn_forward(x, h0, Wx, Wh, b)[0]
361 fWh = lambda Wh: rnn_forward(x, h0, Wx, Wh, b)[0]
362 fb = lambda b: rnn_forward(x, h0, Wx, Wh, b)[0]
363
364 dx_num = compute_numeric_gradient(fx, x, dout)
365 dh0_num = compute_numeric_gradient(fh0, h0, dout)
366 dWx_num = compute_numeric_gradient(fWx, Wx, dout)
367 dWh_num = compute_numeric_gradient(fWh, Wh, dout)
368 db_num = compute_numeric_gradient(fb, b, dout)
369
370 print("dx error: ", rel_error(dx_num, dx))
371 print("dh0 error: ", rel_error(dh0_num, dh0))
372 print("dWx error: ", rel_error(dWx_num, dWx))
373 print("dWh error: ", rel_error(dWh_num, dWh))
374 print("db error: ", rel_error(db_num, db))
375 #%% md
376 ## Vanilla RNN: backward with autograd
377
378 Now we will entirely depend on the PyTorch
    autograd module (`torch.autograd`) to compute the
    backward pass of RNN.
379 `torch.autograd` provides classes and functions
    implementing **automatic differentiation** of
    arbitrary scalar valued functions.
380 It requires minimal changes to the existing code
    - if you pass tensors with `requires_grad=True`
        to the forward function you wrote earlier, you can
        just call ` `.backward(gradient=grad)` on the
        output to compute gradients on the input and
        weights.
381
382 Now we can compare the manual backward pass with
    the autograd backward pass.

```

```
383 Read the code in following cell, and execute it to
    compare your implementation with `torch.autograd
    `.

384 You should get a relative error less than `1e-10`.
385 #%%
386 reset_seed(0)
387
388 N, D, T, H = 2, 3, 10, 5
389
390 # set requires_grad=True
391 x = torch.randn(N, T, D, **to_double,
    requires_grad=True)
392 h0 = torch.randn(N, H, **to_double, requires_grad=
    True)
393 Wx = torch.randn(D, H, **to_double, requires_grad=
    True)
394 Wh = torch.randn(H, H, **to_double, requires_grad=
    True)
395 b = torch.randn(H, **to_double, requires_grad=True
    )
396
397 out, cache = rnn_forward(x, h0, Wx, Wh, b)
398
399 dout = torch.randn(*out.shape, **to_double)
400
401 # Manual backward:
402 with torch.no_grad():
403     dx, dh0, dWx, dWh, db = rnn_backward(dout,
    cache)
404
405 # Backward with autograd: the magic happens here!
406 out.backward(dout)
407
408 dx_auto, dh0_auto, dWx_auto, dWh_auto, db_auto = (
409     x.grad,
410     h0.grad,
411     Wx.grad,
412     Wh.grad,
413     b.grad,
414 )
415
```

```

416 print("dx error: ", rel_error(dx_auto, dx))
417 print("dh0 error: ", rel_error(dh0_auto, dh0))
418 print("dWx error: ", rel_error(dWx_auto, dWx))
419 print("dWh error: ", rel_error(dWh_auto, dWh))
420 print("db error: ", rel_error(db_auto, db))


---


421 #%% md
422 ## RNN Module
423
424 We can now wrap the vanilla RNN implementation
425 into a PyTorch module.
426 Recall from the past assignment/tutorial -- `nn.
427 Module` is a base class for all neural network
428 modules in PyTorch. More details regarding its
429 attributes, functions, and methods could be found
430 [in PyTorch documentation](https://pytorch.org/
431 docs/stable/nn.html?highlight=module#torch.nn.
432 Module).
433
434 In short, the weights and biases are declared in `__init__` and function `forward` will call the `rnn_forward` function from before.
435 The backward function will not be used, and entirely handled by `torch.autograd`.
436 **We have written this part in `RNN` for you but you are highly recommended to go through the code .**

---


437 #%%%
438 from rnn_lstm_captioning import RNN, rnn_forward
439
440
441 N, D, T, H = 2, 3, 10, 5
442
443 x = torch.randn(N, T, D, **to_double)
444 h0 = torch.randn(N, H, **to_double)
445
446 rnn_module = RNN(D, H).to(**to_double)
447
448 # Call forward in module:
449 hn1 = rnn_module(x, h0)
450
451 # Call without module: (but access weights from

```

```
444 module)
445 # Equivalent to above, we won't do this henceforth
446 .
447 Wx, Wh, b = rnn_module.Wx, rnn_module.Wh,
rnn_module.b
448 hn2, _ = rnn_forward(x, h0, Wx, Wh, b)
449 print("Output error with/without module: ",
rel_error(hn1, hn2))
450 #%% md
451 # RNN for image captioning
452
453 You will implement a few necessary tools and
layers in order to build an image captioning model
(class `CaptioningRNN`).
454
455 ## Image Feature Extraction
456
457 The first essential component in an image
captioning model is an encoder that inputs an
image and produces features for decoding the
caption.
458 Here, we use a small [RegNetX-400MF](https://pytorch.org/vision/stable/models.html#torchvision.models.regnet\_x\_400mf) as the backbone so we can
train in reasonable time on Colab. This model is
similar to detector backbone seen in the past
assignment.
459
460 It accepts image batches of shape `(B, C, H, W)`
and outputs spatial features from final layer that
have shape `(B, C, H/32, W/32)` .
461 For vanilla RNN and LSTM, we use the average
pooled features (shape `(B, C)`) for decoding
captions, whereas for attention LSTM we aggregate
the spatial features by learning attention weights
.
462 Checkout the `ImageEncoder` method in `rnn_lstm_captioning.py` to see the initialization
of the model.
463
```

```
464 We use the implementation from torchvision and put
        a very thin wrapper module for our use-case.
465 You do not need to implement anything here – you
        should read and understand the module definition,
        available in `rnn_lstm_captioning.py`

---


466 #%%
467 from rnn_lstm_captioning import ImageEncoder
468
469 model = ImageEncoder(pretrained=True, verbose=True
        ).to(device=DEVICE)

---


470 #%% md
471 ## Word embedding
472 In deep learning systems, we commonly represent
        words using vectors. Each word of the vocabulary
        will be associated with a vector, and these
        vectors will be learned jointly with the rest of
        the system.
473
474 Implement the `WordEmbedding` module to convert
        words (represented by integers) into vectors.
475 Run the following to check your implementation.
        You should see an error on the order of `1e-7` or
        less.

---


476
477 #%%
478 from rnn_lstm_captioning import WordEmbedding
479
480 N, T, V, D = 2, 4, 5, 3
481
482 x = torch.tensor([[0, 3, 1, 2], [2, 1, 0, 3]]).long()
483 W = torch.linspace(0, 1, steps=V * D, **to_double
        ).view(V, D)
484
485 # Copy custom weight vector for sanity check:
486 model_emb = WordEmbedding(V, D).to(**to_double)
487 model_emb.W_embed.data.copy_(W)
488 out = model_emb(x)
489 expected_out = torch.tensor(
490     [
491         [
```

```

492 [0.0, 0.07142857, 0.14285714],
493 [0.64285714, 0.71428571, 0.78571429],
494 [0.21428571, 0.28571429, 0.35714286],
495 [0.42857143, 0.5, 0.57142857],
496 ],
497 [
498 [0.42857143, 0.5, 0.57142857],
499 [0.21428571, 0.28571429, 0.35714286],
500 [0.0, 0.07142857, 0.14285714],
501 [0.64285714, 0.71428571, 0.78571429],
502 ],
503 ],
504 **to_double
505 )
506
507 print("out error: ", rel_error(expected_out, out))
508 #%% md
509 ## Temporal Softmax loss
510
511 In an RNN language model, at every timestep we
512 produce a score for each word in the vocabulary.
512 This score is obtained by applying an affine
513 transform to the hidden state (think `nn.Linear`
514 module).
513 We know the ground-truth word at each timestep, so
515 we use a cross-entropy loss at each timestep.
514 We sum the losses over time and average them over
516 the minibatch.
515
516 However there is one wrinkle: since we operate
517 over minibatches and different captions may have
518 different lengths, we append `<NULL>` tokens to
519 the end of each caption so they all have the same
520 length. We don't want these `<NULL>` tokens to
521 count toward the loss or gradient, so in addition
522 to scores and ground-truth labels our loss
523 function also accepts a `ignore_index` that tells
524 it which index in caption should be ignored when
525 computing the loss.
526
527 Implement the `temporal_softmax_loss` and run the

```

```

518 following cell to check if the implementation is
correct.
519 #%%
520 from rnn_lstm_captioning import
temporal_softmax_loss
521
522
523 def check_loss(N, T, V, p):
524     x = 0.001 * torch.randn(N, T, V)
525     y = torch.randint(V, size=(N, T))
526     mask = torch.rand(N, T)
527     y[mask > p] = 0
528
529     # YOUR_TURN: Implement temporal_softmax_loss
530     print(temporal_softmax_loss(x, y, NULL_index).
item())
531
532
533 check_loss(1000, 1, 10, 1.0) # Should be about 2.
00-2.11
534 check_loss(1000, 10, 10, 1.0) # Should be about
20.6-21.0
535 check_loss(5000, 10, 10, 0.1) # Should be about 2
.00-2.11
536 #%% md
537 ## Captioning Module
538
539 Now we are wrapping everything into the captioning
module. Implement the `CaptioningRNN` module by
following its instructions.
540 This modoule will have a generic structure for RNN
, LST, and attention-based LSTM -- which we
control by providing `cell_type` argument (one of
`["rnn", "lstm", "attn"]`),
541 For now you only need to implement for the case
where `cell_type="rnn"`, you will come back to
this module with other two cases later in this
assignment.
542
543 Also skip the inference function (`CaptioningRNN.
sample`) for now -- only implement `__init__` and

```

```
543 `forward`.  
544 Run the following to check your forward pass using  
      a small test case; you should see difference on  
      the order of `1e-7` or less.  
545 #%%  
546 from rnn_lstm_captioning import CaptioningRNN  
547  
548 reset_seed(0)  
549  
550 N, D, W, H = 10, 400, 30, 40  
551 word_to_idx = {"<NULL>": 0, "cat": 2, "dog": 3}  
552 V = len(word_to_idx)  
553 T = 13  
554  
555 model = CaptioningRNN(  
556     word_to_idx,  
557     input_dim=D,  
558     wordvec_dim=W,  
559     hidden_dim=H,  
560     cell_type="rnn",  
561     ignore_index=NULL_index,  
562 )  
563 # Copy parameters for sanity check:  
564 for k, v in model.named_parameters():  
565     v.data.copy_(torch.linspace(-1.4, 1.3, steps=v  
        .numel()).view(*v.shape))  
566  
567 images = torch.randn(N, 3, *IMAGE_SHAPE)  
568 captions = (torch.arange(N * T) % V).view(N, T)  
569  
570 loss = model(images, captions).item()  
571 expected_loss = 150.6090393066  
572  
573 print("loss: ", loss)  
574 print("expected loss: ", expected_loss)  
575 print("difference: ", rel_error(torch.tensor(loss  
        ), torch.tensor(expected_loss)))  
576 #%% md  
577 ## Overfit small data  
578  
579 To make sure that everything is working as
```

```
579 expected, we can try to overfit this image  
captioning model to a small subset of data.  
580  
581 We have implemented the `train_captioner` function  
which accepts the model and training data, and  
runs a simple training loop - passing data to  
model, collecting training loss, then calling  
`backward()` to obtain gradients. These gradients  
are optimized using the AdamW optimizer.  
582 You can read its implementation in `a2_helper.py  
`.  
583  
584 We will overfit on a subset of 50 examples.  
585 You should see a final loss of less than `0.5` and  
it should be done fairly quickly.  
586 #%%  
587 from a2_helper import train_captioner  
588  
589 reset_seed(0)  
590  
591 # data input  
592 small_num_train = 50  
593 sample_idx = torch.linspace(0, num_train - 1,  
steps=small_num_train).long()  
594 small_image_data = data_dict["train_images"] [  
sample_idx]  
595 small_caption_data = data_dict["trainCaptions"] [  
sample_idx]  
596  
597 # optimization arguments  
598 num_epochs = 80  
599  
600 # create the image captioning model  
601 model = CaptioningRNN(  
602     cell_type="rnn",  
603     word_to_idx=data_dict["vocab"]["token_to_idx"]  
],  
604     input_dim=400, # hard-coded, do not modify  
605     hidden_dim=512,  
606     wordvec_dim=256,  
607     ignore_index=NULL_index,
```

```
608 )
609 model = model.to(**to_float)
610
611 for learning_rate in [1e-3]:
612     print("learning rate is: ", learning_rate)
613     rnn_overfit, _ = train_captioner(
614         model,
615         small_image_data,
616         small_caption_data,
617         num_epochs=num_epochs,
618         batch_size=OVR_BATCH_SIZE,
619         learning_rate=learning_rate,
620         device=DEVICE,
621     )
622 #%% md
623 ## Inference: Sampling Captions
624
625 Unlike classification models, image captioning
models behave very differently at training time
and at test time.
626 At training time, we have access to the ground-
truth caption, so we feed ground-truth words as
input to the RNN at each timestep.
627 At test time, we sample from the distribution over
the vocabulary at each timestep, and feed the
sample as input to the RNN at the next timestep.
628
629 Implement the `CaptioningRNN.sample` for test-time
sampling. After doing so, run the following to
train a captioning model and sample from the model
on both training and validation data.
630
631 ### Train the image captioning model
632
633 Now perform the training on the entire training
set. You should see a final loss less than `2.0`
and each epoch should take ~14s - 44s to run,
depending on the GPU colab assigns you.
634 #%%
635 from a2_helper import train_captioner
636
```

```

637 reset_seed(0)
638
639 # data input
640 small_num_train = num_train
641 sample_idx = torch.randint(num_train, size=(
642     small_num_train,))
643 small_image_data = data_dict["train_images"][
644     sample_idx]
645 small_caption_data = data_dict["train_captions"][
646     sample_idx]
647
648 # create the image captioning model
649 rnn_model = CaptioningRNN(
650     cell_type="rnn",
651     word_to_idx=data_dict["vocab"]["token_to_idx"],
652     input_dim=400, # hard-coded, do not modify
653     hidden_dim=512,
654     wordvec_dim=256,
655     ignore_index=NULL_index,
656 )
657
658 for learning_rate in [1e-3]:
659     print("learning rate is: ", learning_rate)
660     rnn_model_submit, rnn_loss_submit =
661         train_captioner(
662             rnn_model,
663             small_image_data,
664             small_caption_data,
665             num_epochs=60,
666             batch_size=BATCH_SIZE,
667             learning_rate=learning_rate,
668             device=DEVICE,
669             )
670 #####
671 ## Test-time sampling
672 The samples on training data should be very good;
673 the samples on validation data will probably make
674 less sense.
675 #####
676 from a2_helper import decode_captions

```

```

671
672
673 rnn_model.eval()
674
675 for split in ["train", "val"]:
676     sample_idx = torch.randint(
677         0, num_train if split == "train" else
678         num_val, (VIS_BATCH_SIZE,))
679     sample_images = data_dict[split + "_images"][
680         sample_idx]
681     sampleCaptions = data_dict[split + "_captions"]
682         ][sample_idx]
683
684         # decode_captions is loaded from a2_helper.py
685         gt_captions = decode_captions(sample_captions,
686             , data_dict["vocab"]["idx_to_token"])
687
688         generated_captions = rnn_model.sample(
689             sample_images.to(DEVICE))
690         generated_captions = decode_captions(
691             generated_captions, data_dict["vocab"]["idx_to_token"])
692
693         for i in range(VIS_BATCH_SIZE):
694             plt.imshow(sample_images[i].permute(1, 2,
695             0))
696             plt.axis("off")
697             plt.title(
698                 f"[{split}] RNN Generated: {generated_captions[i]}\nGT: {gt_captions[i]}")
699             plt.show()
700
701 %% md
702 # Image Captioning with LSTMs
703
704 So far you have implemented a vanilla RNN and
705 applied it to image captioning.
706 Next we will implement LSTM and use it for image
707 captioning.

```

```

702
703 **LSTM** stands for [Long-Short Term Memory Networks](https://www.researchgate.net/publication/13853244\_Long\_Short-term\_Memory), a variant of vanilla Recurrent Neural Networks.
704 Vanilla RNNs can be tough to train on long sequences due to vanishing and exploding gradients caused by repeated matrix multiplication.
705 LSTMs solve this problem by replacing the simple update rule of the vanilla RNN with a gating mechanism.
706
707 **LSTM Update Rule:** Similar to the vanilla RNN, at each timestep we receive an input  $x_t$  and the previous hidden state  $h_{t-1}$ ; the LSTM also maintains an  $H$ -dimensional *cell state*, so we also receive the previous cell state  $c_{t-1}$ . The learnable parameters of the LSTM are an *input-to-hidden* matrix  $W_x$ , a *hidden-to-hidden* matrix  $W_h$ , a *bias vector*  $b$ .
708
709 At each timestep we first compute an *activation vector*  $a$  as  $a = W_{xx}x_t + W_{hh}h_{t-1} + b$ . We then divide this into four vectors  $a_i, a_f, a_o, a_g$  where  $a_i$  consists of the first  $H$  elements of  $a$ ,  $a_f$  is the next  $H$  elements of  $a$ , etc. We then compute the *input gate*  $g_i$ , *forget gate*  $f$ , *output gate*  $o$  and *block input*  $g$  as
710
711
712 \begin{aligned}
713 i &= \sigma(a_i) \\
714 f &= \sigma(a_f) \\
715 o &= \sigma(a_o) \\
716 g &= \tanh(a_g)
\end{aligned}

```

```

717 \end{align*}
718
719
720 where  $\sigma$  is the sigmoid function and  $\tanh$ 
    is the hyperbolic tangent, both applied
    elementwise.
721
722 Finally we compute the next cell state  $c_t$  and
    next hidden state  $h_t$  as
723
724 $$
725  $c_t = f \odot c_{t-1} + i \odot g$  \hspace{4pc}
726  $h_t = o \odot \tanh(c_t)$ 
727 $$
728
729 where  $\odot$  is the elementwise product of
    vectors.
730
731 In the rest of the notebook we will implement the
    LSTM update rule and apply it to the image
    captioning task.
732 In the code, we assume that data is stored in
    batches so that  $X_t \in \mathbb{R}^{N \times D}$ ,
    and will work with *transposed* versions of the
    parameters:  $W_x \in \mathbb{R}^{D \times 4H}$ ,
 $W_h \in \mathbb{R}^{H \times 4H}$  so that
    activations  $A \in \mathbb{R}^{N \times 4H}$  can be
    computed efficiently as  $A = X_t W_x + H_{t-1}$ 
 $W_h$ 


---


733 #%% md
734 ## LSTM: step forward
735
736 Implement the forward pass for a single timestep
    of an LSTM in the `LSTM.step_forward()` function.
737 This should be similar to the `rnn_step_forward`
    function that you implemented above, but using the
    LSTM update rule instead.
738 Since `LSTM` extends PyTorch `nn.Module`, you don'
    t need to implement backward part!
739
740 Once you are done, run the following to perform a

```

```
740 simple test of your implementation. You should see
    errors on the order of `1e-7` or less.
741 #%%
742 from rnn_lstm_captioning import LSTM
743
744
745 N, D, H = 3, 4, 5
746 x = torch.linspace(-0.4, 1.2, steps=N * D, **  
    to_double).view(N, D)
747 prev_h = torch.linspace(-0.3, 0.7, steps=N * H, **  
    to_double).view(N, H)
748 prev_c = torch.linspace(-0.4, 0.9, steps=N * H, **  
    to_double).view(N, H)
749 Wx = torch.linspace(-2.1, 1.3, steps=4 * D * H, **  
    to_double).view(D, 4 * H)
750 Wh = torch.linspace(-0.7, 2.2, steps=4 * H * H, **  
    to_double).view(H, 4 * H)
751 b = torch.linspace(0.3, 0.7, steps=4 * H, **  
    to_double)
752
753
754 # Create module and copy weight tensors for sanity
    check:
755 model = LSTM(D, H).to(**to_double)
756 model.Wx.data.copy_(Wx)
757 model.Wh.data.copy_(Wh)
758 model.b.data.copy_(b)
759
760 next_h, next_c = model.step_forward(x, prev_h,
    prev_c)
761
762 expected_next_h = torch.tensor(
763     [
764         [0.24635157, 0.28610883, 0.32240467, 0.  
            35525807, 0.38474904],
765         [0.49223563, 0.55611431, 0.61507696, 0.  
            66844003, 0.7159181],
766         [0.56735664, 0.66310127, 0.74419266, 0.  
            80889665, 0.858299],
767     ],
768     **to_double
```

```
769 )
770 expected_next_c = torch.tensor(
771     [
772         [0.32986176, 0.39145139, 0.451556, 0.
773          51014116, 0.56717407],
774         [0.66382255, 0.76674007, 0.87195994, 0.
775          97902709, 1.08751345],
776         [0.74192008, 0.90592151, 1.07717006, 1.
777          25120233, 1.42395676],
778     ],
779     **to_double
780 )
781
782 print("next_h error: ", rel_error(expected_next_h,
783                                     next_h))
784 print("next_c error: ", rel_error(expected_next_c,
785                                     next_c))


---


786 #%%
787 N, D, H, T = 2, 5, 4, 3
788 x = torch.linspace(-0.4, 0.6, steps=N * T * D, **
789                      to_double).view(N, T, D)
790 h0 = torch.linspace(-0.4, 0.8, steps=N * H, **
791                      to_double).view(N, H)
792 Wx = torch.linspace(-0.2, 0.9, steps=4 * D * H, **
793                      to_double).view(D, 4 * H)
794 Wh = torch.linspace(-0.3, 0.6, steps=4 * H * H, **
795                      to_double).view(H, 4 * H)
796 b = torch.linspace(0.2, 0.7, steps=4 * H, **
797                      to_double)
798
799
800 # Create module and copy weight tensors for sanity
801 # check:
```

```
796 model = LSTM(D, H).to(**to_double)
797 model.Wx.data.copy_(Wx)
798 model.Wh.data.copy_(Wh)
799 model.b.data.copy_(b)
800
801 hn = model(x, h0)
802
803 expected_hn = torch.tensor(
804     [
805         [
806             [0.01764008, 0.01823233, 0.01882671, 0
807             .0194232],
808             [0.11287491, 0.12146228, 0.13018446, 0
809             .13902939],
810             [0.31358768, 0.33338627, 0.35304453, 0
811             .37250975],
812             ],
813             [
814                 [0.45767879, 0.4761092, 0.4936887, 0.
815                 51041945],
816                 [0.6704845, 0.69350089, 0.71486014, 0.
817                 7346449],
818                 [0.81733511, 0.83677871, 0.85403753, 0
819                 .86935314],
820                 ],
821             ],
822         ],
823         **to_double
824     )
825
826 print("hn error: ", rel_error(expected_hn, hn))
```

820 %% md

821 ## LSTM captioning model

822

823 Now that you have implemented the `LSTM` module, update the `CaptioningRNN` module (`__init__` and `forward` implementation method ****ONLY****) to also handle the case where `self.cell_type` is `lstm`.

824 ****This should require adding less than 5 lines of code.****

825

826 Once you have done so, run the following to check

```
826 your implementation. You should see a difference  
on the order of `1e-7` or less.  
827 #%%  
828 from rnn_lstm_captioning import CaptioningRNN  
829  
830 N, D, W, H = 10, 400, 30, 40  
831 word_to_idx = {"<NULL>": 0, "cat": 2, "dog": 3}  
832 V = len(word_to_idx)  
833 T = 13  
834  
835 # YOUR_TURN: Implement CaptioningRNN for lstm  
836 model = CaptioningRNN(  
837     word_to_idx,  
838     input_dim=D,  
839     wordvec_dim=W,  
840     hidden_dim=H,  
841     cell_type="lstm",  
842     ignore_index=NULL_index,  
843 )  
844  
845 model = model.to(DEVICE)  
846  
847 for k, v in model.named_parameters():  
848     # print(k, v.shape) # uncomment this to see  
     the weight shape  
849     v.data.copy_(torch.linspace(-1.4, 1.3, steps=v  
     .numel()).view(*v.shape))  
850  
851 images = torch.linspace(  
852     -3.0, 3.0, steps=(N * 3 * IMAGE_SHAPE[0] *  
     IMAGE_SHAPE[1]), **to_float  
853 ).view(N, 3, *IMAGE_SHAPE)  
854 captions = (torch.arange(N * T) % V).view(N, T)  
855  
856 loss = model(images.to(DEVICE), captions.to(DEVICE  
    ))  
857 expected_loss = torch.tensor(146.3161468505)  
858  
859 print("loss: ", loss.item())  
860 print("expected loss: ", expected_loss.item())  
861 print("difference: ", rel_error(loss,
```

```
861 expected_loss))
862 #%% md
863 ## Overfit small data
864 We have written this part for you. Run the
     following to overfit an LSTM captioning model on
     the same small dataset as we used for the RNN
     previously. You should see a final loss less than
     `4` after 80 epochs and it should run fairly
     quickly.
865 #%%
866 from a2_helper import train_captioner
867
868
869 reset_seed(0)
870
871 # Data input.
872 small_num_train = 50
873 sample_idx = torch.linspace(0, num_train - 1,
     steps=small_num_train).long()
874 small_image_data = data_dict["train_images"][
     sample_idx].to(DEVICE)
875 small_caption_data = data_dict["train_captions"][
     sample_idx].to(DEVICE)
876
877 # Create the image captioning model.
878 model = CaptioningRNN(
879     cell_type="lstm",
880     word_to_idx=data_dict["vocab"]["token_to_idx"]
881     ],
882     input_dim=400, # hard-coded, do not modify
883     hidden_dim=512,
884     wordvec_dim=256,
885     ignore_index=NULL_index,
886 )
887 model = model.to(DEVICE)
888 for learning_rate in [1e-2]:
889     print("learning rate is: ", learning_rate)
890     lstm_overfit, _ = train_captioner(
891         model,
892         small_image_data,
```

```
893         small_caption_data,
894         num_epochs=80,
895         batch_size=OVR_BATCH_SIZE,
896         learning_rate=learning_rate,
897     )
898 #%% md
899 ## Caption sampling
900
901 Modify the `CaptioningRNN.sample` method in class
    to handle the case where `self.cell_type` is `lstm`. **This should take fewer than 10 lines of
    code.**
902 When you are done, run the following cells to
    train the captioning model first, then sample some
    captions from your model during test time.
903
904 ### Train the net
905
906 Perform the training on the entire training set.
    You should see a final loss less than `1.8`. Each
    epoch should take ~7s - 14s to run, depending on
    the GPU
907 colab assigns you.
908 #%%
909 reset_seed(0)
910
911 # data input
912 small_num_train = num_train
913 sample_idx = torch.randint(num_train, size=(
    small_num_train,))
914 small_image_data = data_dict["train_images"][
    sample_idx]
915 small_caption_data = data_dict["trainCaptions"][
    sample_idx]
916
917 # create the image captioning model
918 lstm_model = CaptioningRNN(
919     cell_type="lstm",
920     word_to_idx=data_dict["vocab"]["token_to_idx"]
921     ],
921     input_dim=400, # hard-coded, do not modify
```

```
922     hidden_dim=512,  
923     wordvec_dim=256,  
924     ignore_index=NULL_index,  
925 )  
926 lstm_model = lstm_model.to(DEVICE)  
927  
928 for learning_rate in [1e-3]:  
929     print("learning rate is: ", learning_rate)  
930     lstm_model_submit, lstm_loss_submit =  
931         train_captioner(  
932             lstm_model,  
933             small_image_data,  
934             small_caption_data,  
935             num_epochs=60,  
936             batch_size=BATCH_SIZE,  
937             learning_rate=learning_rate,  
938             device=DEVICE,  
939         )  
940 #%% md  
940 ### Test-time sampling  
941 As with the RNN, the samples on training data  
should be very good; the samples on validation  
data will probably make less sense.  
942 #%%  
943 from a2_helper import decode_captions  
944  
945  
946 lstm_model.eval()  
947  
948 for split in ["train", "val"]:  
949     sample_idx = torch.randint(  
950         0, num_train if split == "train" else  
951         num_val, (VIS_BATCH_SIZE,) )  
952     sample_images = data_dict[split + "_images"][sample_idx]  
953     sample_captions = data_dict[split + "_captions"] [sample_idx]  
954  
955     # decode_captions is loaded from a2_helper.py  
956     gt_captions = decode_captions(sample_captions)
```

```

956 , data_dict["vocab"]["idx_to_token"])
957     lstm_model.eval()
958     generatedCaptions = lstm_model.sample(
959         sample_images.to(DEVICE))
960     generatedCaptions = decodeCaptions(
961         generatedCaptions, data_dict["vocab"]["idx_to_token"])
962
963     for i in range(VIS_BATCH_SIZE):
964         plt.imshow(sample_images[i].permute(1, 2, 0))
965         plt.axis("off")
966         plt.title(
967             f"[{split}] LSTM Generated: {generatedCaptions[i]}\nGT: {gtCaptions[i]}")
968
969     plt.show()
970 #%% md
971 # Attention LSTM
972 Attention LSTM essentially adds an attention input
$ x_{attn}^t \in \mathbb{R}^H $ into LSTM, along with
$ x_t \in \mathbb{R}^D $ and the previous hidden
state $ h_{t-1} \in \mathbb{R}^H $.
973
974 To get the attention input $ x_{attn}^t $, here we
adopt a method called `scaled dot-product
attention`, as covered in the lecture. We first
project the CNN feature activation from $ \mathbb{R}^{400 \times 4 \times 4} $ to $ \mathbb{R}^{H \times 4 \times 4} $
using an affine layer. Given the
projected activation $ A \in \mathbb{R}^{H \times 4 \times 4} $ and the LSTM hidden state from the
previous time step $ h_{t-1} $, we formulate the
attention weights on $ A $ at time step $ t $ as $ M_{attn}^t = h_{t-1} A / \sqrt{H} \in \mathbb{R}^{4 \times 4} $.
975
976 To simplify the formulation here, we flatten the
spatial dimensions of $ A $ and $ M_{attn}^t $ which
gives $ \tilde{A} \in \mathbb{R}^{H \times 16} $ and $ \tilde{M}_{attn}^t \in \mathbb{R}^{4 \times 16} $.

```

```

976  $\tilde{M}^t$ _{attn} =  $h_{t-1}A \in \mathbb{R}^{16}$ .
977 We add a **`softmax`** activation function on  $\tilde{M}^t$ _{attn} so that the attention weights at each time step are normalized and sum up to one.
978
979 The attention embedding given the attention weights is then  $x_{attn}^t = \tilde{A}\tilde{M}^t$ _{attn}  $\in \mathbb{R}^H$ . Next, you will implement a batch version of the attention layer we have described here.


---


980 #%% md
981 ## Scaled dot-product attention
982 Implement the `dot_product_attention` function. Given the LSTM hidden state from the previous time step `prev_h` (or  $h_{t-1}$ ) and the projected CNN feature activation `A`, compute the attention weights `attn_weights` (or  $\tilde{M}^t$ _{attn}) with a reshaping to  $\mathbb{R}^{4 \times 4}$  attention embedding output `attn` (or  $x_{attn}^t$ ) using the formulation we provided.
983
984 When you are done, run the following to check your implementation. You should see an error on the order of `1e-7` or less.


---


985 #%%
986 from rnn_lstm_captioning import
    dot_product_attention
987
988
989 N, H = 2, 5
990 D_a = 4
991
992 prev_h = torch.linspace(-0.4, 0.6, steps=N * H
    , **to_double).view(N, H)
993 A = torch.linspace(-0.4, 1.8, steps=N * H * D_a
    * D_a, **to_double).view(
994     N, H, D_a, D_a
995 )
996
997 # YOUR_TURN: Implement dot_product_attention

```

```
998 attn, attn_weights = dot_product_attention(prev_h  
    , A)  
999  
1000 expected_attn = torch.tensor(  
1001     [  
1002         [-0.29784344, -0.07645979, 0.14492386, 0.  
            36630751, 0.58769115],  
1003         [0.81412643, 1.03551008, 1.25689373, 1.  
            47827738, 1.69966103],  
1004     ],  
1005     **to_double  
1006 )  
1007 expected_attn_weights = torch.tensor(  
1008     [  
1009         [  
1010             [0.06511126, 0.06475411, 0.06439892,  
                0.06404568],  
1011             [0.06369438, 0.06334500, 0.06299754,  
                0.06265198],  
1012             [0.06230832, 0.06196655, 0.06162665,  
                0.06128861],  
1013             [0.06095243, 0.06061809, 0.06028559,  
                0.05995491],  
1014         ],  
1015         [  
1016             [0.05717142, 0.05784357, 0.05852362,  
                0.05921167],  
1017             [0.05990781, 0.06061213, 0.06132473,  
                0.06204571],  
1018             [0.06277517, 0.06351320, 0.06425991,  
                0.06501540],  
1019             [0.06577977, 0.06655312, 0.06733557,  
                0.06812722],  
1020         ],  
1021     ],  
1022     **to_double  
1023 )  
1024  
1025 print("attn error: ", rel_error(expected_attn,  
    attn))  
1026 print("attn_weights error: ", rel_error(
```

```

1026 expected_attn_weights, attn_weights))
1027 #%% md
1028 ## Attention LSTM: step forward
1029
1030 Implement `AttentionLSTM.step_forward()` by
    following its instructions and input
    specifications.
1031 It is mostly similar to `LSTM.step_forward()` but
    has extra attention input `attn` (or  $x_{\text{attn}}$ )
    and its embedding weight matrix `Wattn` (or  $W_{\text{attn}}$ ),
1032 these are defined in `AttentionLSTM.__init__()` .
1033 Hence, at each timestep the *activation vector*
 $\mathbf{a} \in \mathbb{R}^{4H}$  in LSTM cell is formulated
    as:
1034
1035  $\mathbf{a} = W_{xx} \mathbf{x}_t + W_{hh} \mathbf{h}_{t-1} + W_{\text{attn}} \mathbf{x}_{\text{attn}}^t + b$ .
1036
1037
1038 **This should require adding less than 5 lines of
    code.**
1039 Once you are done, run the following to perform a
    simple test of your implementation. You should
    see errors on the order of `1e-8` or less.
1040 #%%
1041 from rnn_lstm_captioning import AttentionLSTM
1042
1043
1044 N, D, H = 3, 4, 5
1045
1046 x = torch.linspace(-0.4, 1.2, steps=N * D, **
    to_double).view(N, D)
1047 prev_h = torch.linspace(-0.3, 0.7, steps=N * H
    , **to_double).view(N, H)
1048 prev_c = torch.linspace(-0.4, 0.9, steps=N * H
    , **to_double).view(N, H)
1049 attn = torch.linspace(0.6, 1.8, steps=N * H, **
    to_double).view(N, H)
1050
1051 Wx = torch.linspace(-2.1, 1.3, steps=4 * D * H
    , **to_double).view(D, 4 * H)

```

```
1052 Wh = torch.linspace(-0.7, 2.2, steps=4 * H * H  
    , **to_double).view(H, 4 * H)  
1053 b = torch.linspace(0.3, 0.7, steps=4 * H, **  
    to_double)  
1054 Wattn = torch.linspace(1.3, 4.2, steps=4 * H * H  
    , **to_double).view(H, 4 * H)  
1055  
1056 # Create module and copy weight tensors for  
    sanity check:  
1057 model = AttentionLSTM(D, H).to(**to_double)  
1058 model.Wx.data.copy_(Wx)  
1059 model.Wh.data.copy_(Wh)  
1060 model.b.data.copy_(b)  
1061 model.Wattn.data.copy_(Wattn)  
1062  
1063 next_h, next_c = model.step_forward(x, prev_h,  
    prev_c, attn)  
1064  
1065  
1066 expected_next_h = torch.tensor(  
1067     [  
1068         [0.53704256, 0.59980774, 0.65596820, 0.  
        70569729, 0.74932626],  
1069         [0.78729857, 0.82010653, 0.84828362, 0.  
        87235677, 0.89283167],  
1070         [0.91017981, 0.92483119, 0.93717126, 0.  
        94754073, 0.95623746],  
1071     ],  
1072     **to_double  
1073 )  
1074 expected_next_c = torch.tensor(  
1075     [  
1076         [0.59999328, 0.69285041, 0.78570758, 0.  
        87856479, 0.97142202],  
1077         [1.06428558, 1.15714276, 1.24999992, 1.  
        34285708, 1.43571424],  
1078         [1.52857143, 1.62142857, 1.71428571, 1.  
        80714286, 1.90000000],  
1079     ],  
1080     **to_double  
1081 )
```

```
1082
1083 print("next_h error: ", rel_error(expected_next_h
1084 , next_h))
1084 print("next_c error: ", rel_error(expected_next_c
1085 , next_c))
1085 #%% md
1086 ## Attention LSTM: forward
1087
1088 Now, implement the `AttentionLSTM.forward()` function to run an attention-based LSTM on an entire timeseries of data.
1089 You will have to use the `dot_product_attention` function from outside this module.
1090
1091 When you are done, run the following to check your implementation. You should see an error on the order of `1e-8` or less.
1092 #%%
1093 N, D, H, T = 2, 5, 4, 3
1094 D_a = 4
1095
1096 x = torch.linspace(-0.4, 0.6, steps=N * T * D, **to_double).view(N, T, D)
1097 A = torch.linspace(-0.4, 1.8, steps=N * H * D_a
1098 * D_a, **to_double).view(
1099 N, H, D_a, D_a
1099 )
1100
1101 Wx = torch.linspace(-0.2, 0.9, steps=4 * D * H
1102 , **to_double).view(D, 4 * H)
1102 Wh = torch.linspace(-0.3, 0.6, steps=4 * H * H
1103 , **to_double).view(H, 4 * H)
1103 Wattn = torch.linspace(1.3, 4.2, steps=4 * H * H
1104 , **to_double).view(H, 4 * H)
1104 b = torch.linspace(0.2, 0.7, steps=4 * H, **
1105 to_double)
1105
1106
1107 # Create module and copy weight tensors for
1108 # sanity check:
1108 model = AttentionLSTM(D, H).to(**to_double)
```

```

1109 model.Wx.data.copy_(Wx)
1110 model.Wh.data.copy_(Wh)
1111 model.b.data.copy_(b)
1112 model.Wattn.data.copy_(Wattn)
1113
1114 # YOUR_TURN: Implement attention_forward
1115 hn = model(x, A)
1116
1117 expected_hn = torch.tensor(
1118     [
1119         [
1120             [0.56141729, 0.70274849, 0.80000386,
1121              0.86349400],
1122              [0.89556391, 0.92856726, 0.94950579,
1123              0.96281018],
1124              [0.96792077, 0.97535465, 0.98039623,
1125              0.98392994],
1126              ],
1127              [
1128                  [0.95065880, 0.97135490, 0.98344373,
1129                  0.99045552],
1130                  [0.99317679, 0.99607466, 0.99774317,
1131                  0.99870293],
1132                  [0.99907382, 0.99946784, 0.99969426,
1133                  0.99982435],
1134                  ],
1135                  ],
1136                  **to_double
1137 )
1138
1139 print("h error: ", rel_error(expected_hn, hn))
1140 %% md
1141 ## Attention LSTM captioning model
1142
1143 With all your implementation done so far, you can
1144 finally update the implementation of
1145 `CaptioningRNN.__init__` and `CaptioningRNN.
1146 forward` methods once again.
1147 This time, handle the case where `self.cell_type
1148 ` is `attn`. **This should require adding less
1149 than 10 lines of code.**
```

```
1139
1140 Once you have done so, run the following to check
    your implementation. You should see a difference
    on the order of `1e-7` or less.
1141 #%%
1142 from rnn_lstm_captioning import CaptioningRNN
1143
1144
1145 reset_seed(0)
1146
1147 N, D, W, H = 10, 400, 30, 40
1148 word_to_idx = {"<NULL>": 0, "cat": 2, "dog": 3}
1149 V = len(word_to_idx)
1150 T = 13
1151
1152 # YOUR_TURN: Modify CaptioningRNN for attention
1153 model = CaptioningRNN(
1154     word_to_idx,
1155     input_dim=D,
1156     wordvec_dim=W,
1157     hidden_dim=H,
1158     cell_type="attn",
1159     ignore_index=NULL_index,
1160 )
1161 model = model.to(DEVICE)
1162
1163 for k, v in model.named_parameters():
1164     # print(k, v.shape) # uncomment this to see
        the weight shape
1165     v.data.copy_(torch.linspace(-1.4, 1.3, steps=v.numel()).view(*v.shape))
1166
1167 images = torch.linspace(
1168     -3.0, 3.0, steps=(N * 3 * IMAGE_SHAPE[0] *
1169     IMAGE_SHAPE[1]))
1170 captions = (torch.arange(N * T) % V).view(N, T)
1171
1172 loss = model(images.to(DEVICE), captions.to(
1173     DEVICE))
1173 expected_loss = torch.tensor(8.0156393051)
```

```
1174
1175 print("loss: ", loss.item())
1176 print("expected loss: ", expected_loss.item())
1177 print("difference: ", rel_error(loss,
1178     expected_loss))
1178 %% md
1179 ## Overfit small data
1180 We have written this part for you. Run the
1181   following to overfit an Attention LSTM captioning
1182   model on the same small dataset as we used for
1183   the RNN previously. You should see a final loss
1184   less than `9`.
1185
1186 # data input
1187 small_num_train = 50
1188 sample_idx = torch.linspace(0, num_train - 1,
1189   steps=small_num_train).long()
1190 small_image_data = data_dict["train_images"][
1191   sample_idx]
1192 small_caption_data = data_dict["trainCaptions"][
1193   sample_idx]
1194
1195 # create the image captioning model
1196 model = CaptioningRNN(
1197   cell_type="attn",
1198   word_to_idx=data_dict["vocab"]["token_to_idx"],
1199   input_dim=400, # hard-coded, do not modify
1200   hidden_dim=512,
1201   wordvec_dim=256,
1202   ignore_index=NULL_index,
1203 )
1204
1205 for learning_rate in [1e-3]:
1206   print("learning rate is: ", learning_rate)
1207   attn_overfit, _ = train_captioner(
```

```
1206         model,
1207         small_image_data,
1208         small_caption_data,
1209         num_epochs=80,
1210         batch_size=OVR_BATCH_SIZE,
1211         learning_rate=learning_rate,
1212         device=DEVICE,
1213     )
1214 #%% md
1215 ## Caption sampling
1216
1217 Modify the `CaptioningRNN.sample` method to
    handle the case where `self.cell_type` is `attn
    `. **This should take fewer than 10 lines of code
    .**
1218
1219 When you are done run the following to train a
    captioning model and sample from the model on
    some training and validation set samples.
1220
1221 ### Train the net
1222
1223 Now, perform the training on the entire training
    set. You should see a final loss less than `0.5
    `. Each epoch should take ~8s to run.
1224 #%%
1225 reset_seed(0)
1226
1227 # data input
1228 small_num_train = num_train
1229 sample_idx = torch.randint(num_train, size=(
    small_num_train,))
1230 small_image_data = data_dict["train_images"][
    sample_idx]
1231 small_caption_data = data_dict["trainCaptions"][
    sample_idx]
1232
1233 # create the image captioning model
1234 attn_model = CaptioningRNN(
1235     cell_type="attn",
1236     word_to_idx=data_dict["vocab"]["token_to_idx"]
```

```

1236 ],
1237     input_dim=400, # hard-coded, do not modify
1238     hidden_dim=512,
1239     wordvec_dim=256,
1240     ignore_index=NULL_index,
1241 )
1242 attn_model = attn_model.to(DEVICE)
1243
1244 for learning_rate in [1e-3]:
1245     print("learning rate is: ", learning_rate)
1246     attn_model_submit, attn_loss_submit =
1247         train_captioner(
1248             attn_model,
1249             small_image_data,
1250             small_caption_data,
1251             num_epochs=60,
1252             batch_size=BATCH_SIZE,
1253             learning_rate=learning_rate,
1254             device=DEVICE,
1255         )
1256 %% md
1257 ### Test-time sampling and visualization
1258 As with RNN and LSTM, the samples on training
1259 data should be very good; the samples on
1260 validation data will probably make less sense.
1261 We use the `attention_visualizer` function from `dl/utils.py` to visualize the attended regions
1262 per generated word. Note that sometimes the
1263 attended regions (brighter) might not make much
1264 sense partially due to our low resolution image
1265 input. In real applications, the attended
1266 regions are more accurate.
1267 %% 
1268 # Sample a minibatch and show the reshaped
1269 # 112x112 images,
1270 # GT captions, and generated captions by your
1271 # model.
1272
1273
1274 from torchvision import transforms
1275 from torchvision.utils import make_grid

```

```

1266
1267 for split in ["train", "val"]:
1268     sample_idx = torch.randint(
1269         0, num_train if split == "train" else
1270         num_val, (VIS_BATCH_SIZE,))
1271     sample_images = data_dict[split + "_images"][
1272         sample_idx]
1273     sampleCaptions = data_dict[split + "
1274         _captions"][sample_idx]
1275         # decode_captions is loaded from a2_helper.py
1276         gt_captions = decode_captions(sampleCaptions
1277         , data_dict["vocab"]["idx_to_token"])
1278         attn_model.eval()
1279         generated_captions, attn_weights_all =
1280             attn_model.sample(sample_images.to(DEVICE))
1281         generated_captions = decode_captions(
1282             generated_captions, data_dict["vocab"]["
1283             idx_to_token"])
1284         )
1285
1286         for i in range(VIS_BATCH_SIZE):
1287             plt.imshow(sample_images[i].permute(1, 2
1288             , 0))
1289             plt.axis("off")
1290             plt.title(
1291                 "%s\nAttention LSTM Generated:%s\nGT
1292                 :%s"
1293                 % (split, generated_captions[i],
1294                     gt_captions[i]))
1295             )
1296             plt.show()
1297
1298         tokens = generated_captions[i].split(" ")
1299
1300         vis_attn = []
1301         for j in range(len(tokens)):
1302             img = sample_images[i]
1303             attn_weights = attn_weights_all[i][j]
1304             token = tokens[j]

```

```
1298         img_copy = attention_visualizer(img,
1299                                         attn_weights, token)
1300                                         vis_attn.append(transforms.ToTensor()
1301                                         ((img_copy)))
1300
1301         plt.rcParams["figure.figsize"] = (20.0,
1302                                         20.0)
1302         vis_attn = make_grid(vis_attn, nrow=8)
1303         plt.imshow(torch.flip(vis_attn, dims=(0
1304                                         ,)).permute(1, 2, 0))
1304         plt.axis("off")
1305         plt.show()
1306         plt.rcParams["figure.figsize"] = (10.0, 8
1307                                         .0)
1307 #%% md
1308 # Save results for submission
1309
1310 Once you have finished all your implementation,
    run "Runtime -> Restart and run all..." to re-run
        all cells and display outputs.
1311 Make sure all outputs are displayed properly and
    the outputs are same as what you expected!
1312
1313 Once all the cells are completed, execute the
    following cell to save the final losses for
        submission.
1314 #%%
1315 submission = {
1316     "rnn_losses": rnn_loss_submit,
1317     "lstm_losses": lstm_loss_submit,
1318     "attn_losses": attn_loss_submit,
1319 }
1320
1321 output_dir = os.path.abspath("./submissions")
1322 os.makedirs(output_dir, exist_ok=True)
1323
1324 submission_path = os.path.join(output_dir, "
1325     rnn_lstm_attention_submission.pt")
1325 torch.save(submission, submission_path)
1326
1327 print("Saved to:", submission_path)
```

```
1 #%% md
2 # Deep Learning Assignment 2 Part 2: Transformers
3 #%% md
4 ### Transformers ([Attention is all you need](https://arxiv.org/pdf/1706.03762.pdf))
5
6 To this point we have seen RNNs, which excel at sequence to sequence task but have two major drawbacks.
7 First, they can suffer from vanishing gradients for long sequences.
8 Second, they can take a long time to train due to sequential dependencies between hidden states which does not take advantage of the massively parallel architecture of modern GPUs.
9 The first issue is largely addressed by alternate RNN architectures (LSTMs, GRUs) but not the second.
10
11 Transformers solve these problems up to a certain extent by enabling to process the input parallelly during training with long sequences. Though the computation is quadratic with respect to the input sequence length, it still managable with modern GPUs.
12
13 In this notebook, we will implement Transformers model step-by-step by referencing the original paper, [Attention is all you need](https://arxiv.org/pdf/1706.03762.pdf). We will also use a toy dataset to solve a vector-to-vector problem which is a subset of sequence-to-sequence problem.
14 #%% md
15 ## Table of Contents
16
17 This assignment has 4 parts. In the class we learned about Encoder based Transformers but often we use an Encoder and a Decoder for sequence to sequence task. In this notebook, you will learn how to implement an Encoder-Decoder based Transformers in a step-by-step manner. We will implement a simpler version here, where the simplicity arise
```

```
17 from the task that we are solving, which is a
    vector-to-vector task. This essentially means that
    the length of input and output sequence is **fixed**
    and we dont have to worry about variable length
    of sequences. This makes the implementation simpler
    .
18
19 1. **Part I (Preparation)**: We will preprocess a
    toy dataset that consists of input arithmetic
    expression and an output result of the expression
20 1. **Part II (Implement Transformer blocks)**: we
    will look how to implement building blocks of a
    Transformer. It will consist of following blocks
21     1. MultiHeadAttention
22     2. FeedForward
23     3. LayerNorm
24     4. Encoder Block
25     5. Decoder Block
26 1. **Part III (Data Loading)**: We will use the
    preprocessing functions in part I and the
    positional encoding module to construct the
    Dataloader.
27 1. **Part IV (Train a model)**: In the last part we
    will look at how to fit the implemented
    Transformer model to the toy dataset.
28
29 You can run all things on CPU till part 3. Part 4
    requires GPU and while changing the runtime for
    this part, you would also have to run all the
    previous parts as part 4 has dependency on previous
    parts.


---


30 #%% md
31 # Part I. Preparation


---


32 #%% md
33 Before getting started we need to run some
    boilerplate code to set up our environment. You'll
    need to rerun this setup code each time you start
    the notebook.
34
35 First, run this cell load the [autoreload](https://ipython.readthedocs.io/en/stable/config/extensions/
```

```
35 autoreload.html?highlight=autoreload) extension.  
This allows us to edit `*.py` source files, and re-import them into the notebook for a seamless editing and debugging experience.

---



```
36 #%%
37 # load_ext autoreload
38 # autoreload 2

```
39 #%% md  
40 ### Google Colab Setup  
41  
42 Next we need to run a few commands to set up our environment on Google Colab. If you are running this notebook on a local machine you can skip this section.  
43  
44 Run the following cell to mount your Google Drive. Follow the link, sign in to your Google account (the same account you used to store this notebook!) and copy the authorization code into the text box that appears below.

---



```
45 #%%
46 # Colab-only setup (Google Drive mount). Safe to ignore locally.
47 try:
48 from google.colab import drive # type: ignore
49 drive.mount("/content/drive")
50 except Exception:
51 print("Not running on Colab: skipping Google Drive mount")

```
52 #%% md  
53 Now recall the path in your Google Drive where you uploaded this notebook, fill it in below. If everything is working correctly then running the following cell should print the filenames from the assignment:  
54  
55 ````  
56 ["dl", "a2_helper.py", "rnn_lstm_captioning.ipynb", "rnn_lstm_captioning.py", "transformers.py", "Transformers.ipynb", "two_digit_op.json"]  
57 ````
```


```


```


```


```

```
58 #%%
59 import sys
60 from pathlib import Path
61
62 # Local (non-Colab) setup: assume this notebook is
63 # inside the assignment folder.
63 ASSIGNMENT_DIR = Path.cwd().resolve()
64
65 # If launched from elsewhere, try to find the
66 # folder that contains transformers.py
66 if not (ASSIGNMENT_DIR / "transformers.py").exists():
67     for candidate in [
68         ASSIGNMENT_DIR / "assignment2",
69         ASSIGNMENT_DIR.parent,
70         ASSIGNMENT_DIR.parent / "assignment2",
71     ]:
72         if (candidate / "transformers.py").exists():
73             ASSIGNMENT_DIR = candidate.resolve()
74             break
75
76 print("Using ASSIGNMENT_DIR:", ASSIGNMENT_DIR)
77 sys.path.append(str(ASSIGNMENT_DIR))
78 #%% md
79 Once you have successfully mounted your Google
Drive and located the path to this assignment, run
the following cell to allow us to import from the `.
py` files of this assignment. If it works correctly
, it should print the message:
80
81 ```
82 Hello from Transformers.py!
83 ```
84
85 as well as the last edit time for the file `.
Transformers.py` .
86 #%%
87 import os
88 import time
89 from transformers import hello_transformers
```

```
90
91
92 os.environ["TZ"] = "Asia/Jerusalem"
93 time.tzset()
94 hello_transformers()
95
96 transformers_path = str(ASSIGNMENT_DIR / "
97     transformers.py")
98 print("transformers.py last edited on %s" %
99     transformers_edit_time)
100 #%%
101 import torch
102 import torch.nn.functional as F
103 from torch import Tensor
104 from torch import nn
105 import torch
106
107 from torch import nn
108 import torch.nn.functional as F
109
110 from dl.utils import (
111     reset_seed,
112     tensor_to_image,
113     attention_visualizer,
114 )
115 from dl.grad import rel_error,
116     compute_numeric_gradient
117 import matplotlib.pyplot as plt
118 import time
119 from IPython.display import Image
120
121 # for plotting
122 %matplotlib inline
123 plt.rcParams["figure.figsize"] = (10.0, 8.0)  #
124     set default size of plots
125 plt.rcParams["image.interpolation"] = "nearest"
126 plt.rcParams["image.cmap"] = "gray"
```

```

126 %% md
127 We will use the GPU to accelerate our computation
    . Run this cell to make sure you are using a GPU.
128
129 We will be using `torch.float = torch.float32` for
    data and `torch.long = torch.int64` for labels.
130
131 Please refer to https://pytorch.org/docs/stable/tensor\_attributes.html#torch-dtype for more
details about data types.
132 %%
133
134 # Device selection logic (unchanged)
135 if torch.backends.mps.is_available() and torch.
    backends.mps.is_built():
136     DEVICE = torch.device("mps")
137     print("Using Apple Silicon GPU (MPS)")
138 elif torch.cuda.is_available():
139     DEVICE = torch.device("cuda")
140     print("Using CUDA GPU")
141 else:
142     DEVICE = torch.device("cpu")
143     print("Using CPU")
144
145 to_float = {"dtype": torch.float32, "device":
    DEVICE}
146
147 DOUBLE_DEVICE = DEVICE
148 if DEVICE.type == "mps":
149     DOUBLE_DEVICE = torch.device("cpu")
150
151 to_double = {"dtype": torch.float64, "device":
    DOUBLE_DEVICE}
152 %% md
153 ### Load the toy data
154 %% md
155 As Transformers perform very well on sequence to
    sequence task, we will implement it on a toy task
    of Arithmetic operations. We will use transformer
    models to perform addition and subtraction of two
    integers, where the absolute value of an integer

```

```
155 is at most 50. A simple example is to perform the
computation `‐5 + 2` using a Transformer model and
getting the corect result as `‐3`. As there can
be multiple ways to solve this problem, we will
see how we can pose this as a sequence to sequence
problem and solve it using Transformers model.
Note that we had to reduce the complexity of the
problem to make the Transformer work within the
constrained resources of Colab.
```

```
156
```

```
157 Lets take a look at the data first:
```

```
158 #%%
```

```
159 from a2_helper import get_toy_data
```

```
160
```

```
161 # load the data using helper function
```

```
162 data = get_toy_data(str(ASSIGNMENT_DIR / "
two_digit_op.json"))
```

```
163 #%% md
```

```
164 ### Looking at the first four examples
```

```
165
```

```
166 Below are the first four samples in the data
```

```
167 #%%
```

```
168 num_examples = 4
```

```
169 for q, a in zip(
```

```
170     data["inp_expression"][:num_examples],
171     data["out_expression"][:num_examples]
```

```
172 ):
```

```
173     print("Expression: " + q + " Output: " + a)
```

```
174 #%% md
```

```
175 ## What do these examples mean:
```

```
176
```

```
177 Lets look at first and third examples here and
understand what they represent:
```

```
178
```

```
179 - Expression: `BOS NEGATIVE 30 subtract NEGATIVE
34 EOS` Output: `BOS POSITIVE 04 EOS`: The
expression here is $(-30) - (-34)$. There are two
notions of the symbol `+` here: one is to denote
the sign of the number and other is the operation
of addition between two integers. To simplify the
problem for the neural network, we have denoted
```

179 them with different text tokens. The output of \$(-30) - (-34)\$ is \$+4\$. Here `BOS` and `EOS` refer to begining of sequence and end of sequence

180 - Similarly, the second expression, `BOS NEGATIVE 34 add NEGATIVE 15 EOS` Output: `BOS NEGATIVE 49 EOS` means that we are doing the computation as \$(-34) + (-15)\$. As above, the symbol `-` here represents two things: first is the sign of an integer and second is the operation between two integers. Again, we have represented with different tokens to simplify the problem for the neural network. The output here is -49. Here `BOS` and `EOS` refer to begining of sequence and end of sequence

181

182 Now that we have a grasp on what is the data, lets head to preprocess the data, as the neural networks don't really understand strings, we need to represent them as numbers.

183

184 ## Pre-processing the data

185 We need to convert the raw input sequence into a format that can be processed with a neural network .

186 Concretely, we need to convert a human-readable string (e.g. `BOS NEGATIVE 30 subtract NEGATIVE 34 EOS`) into a sequence of **tokens**, each of which will be an integer.

187 The process of converting an input string into a sequence of tokens is known as **tokenization**.

188

189 Before we can tokenize any particular sequence, we first need to build a **vocabulary**;

190 this is an exhaustive list of all tokens that appear in our dataset, and a mapping from each token to a unique integer value.

191 In our case, our vocabulary with consist of 16 elements: one entry for each digit `0` to `9`, two tokens to represent the sign of a number (`POSITIVE` and `NEGATIVE`), two tokens representing the addition and subtraction operations (`add` ,

191 and `subtract`), and finally two special tokens representing the start and end of the sequence (`BOS`, `EOS`).

192

193 We typically represent the vocabulary with a pair of data structures.

194 First is a list of all the string tokens (`vocab` below), such that `vocab[i] = s` means that the string `s` has been assigned the integer value `i`. This allows us to look up the string associated with any numeric index `i`.

195 We also need a data structure that enables us to map in the other direction: given a string `s`, find the index `i` to which it has been assigned. This is typically represented as a hash map (`dict` object in Python) whose keys are strings and whose values are the indices assigned to those strings.

196 You will implement the function `generate_token_dict` that inputs the list `vocab` and returns a dict `convert_str_to_token` giving this mapping.

197

198 Once you have built the vocab, then you can implement the function `preprocess_input_sequence` which uses the vocab data structures to convert an input string into a list of integer tokens.

199 #%%

200 # Create vocab

201 SPECIAL_TOKENS = ["POSITIVE", "NEGATIVE", "add", "subtract", "BOS", "EOS"]

202 vocab = ["0", "1", "2", "3", "4", "5", "6", "7", "8", "9"] + SPECIAL_TOKENS

203 #%% md

204 To generate the hash map and then process the input string using them, complete the `generate_token_dict`, `preprocess_input_sequence` functions in the python files for this exercise:

205

206 You should see exact zero errors here

207 #%%

```
208 from transformers import generate_token_dict
209
210 convert_str_to_tokens = generate_token_dict(vocab)
211
212 try:
213     assert convert_str_to_tokens["0"] == 0
214 except:
215     print("The first element does not map to 0.
Please check the implementation")
216
217 try:
218     assert convert_str_to_tokens["EOS"] == 15
219 except:
220     print("The last element does not map to 2004.
Please check the implementation")
221
222 print("Dictionary created successfully!")
223 #%%
224 from transformers import prepocess_input_sequence
225
226 convert_str_to_tokens = generate_token_dict(vocab)
227
228 ex1_in = "BOS POSITIVE 0333 add POSITIVE 0696 EOS"
229 ex2_in = "BOS POSITIVE 0673 add POSITIVE 0675 EOS"
230 ex3_in = "BOS NEGATIVE 0286 subtract NEGATIVE 0044
EOS"
231 ex4_in = "BOS NEGATIVE 0420 add POSITIVE 0342 EOS"
232
233 ex1_out = "BOS POSITIVE 1029 EOS"
234 ex2_out = "BOS POSITIVE 1348 EOS"
235 ex3_out = "BOS NEGATIVE 0242 EOS"
236 ex4_out = "BOS NEGATIVE 0078 EOS"
237
238 ex1_inp_preprocessed = torch.tensor(
239     prepocess_input_sequence(ex1_in,
240     convert_str_to_tokens, SPECIAL_TOKENS)
241 )
242 ex2_inp_preprocessed = torch.tensor(
243     prepocess_input_sequence(ex2_in,
244     convert_str_to_tokens, SPECIAL_TOKENS)
245 )
```

```
244 ex3_inp_preprocessed = torch.tensor(  
245     preprocess_input_sequence(ex3_in,  
246         convert_str_to_tokens, SPECIAL_TOKENS)  
247 )  
248 ex4_inp_preprocessed = torch.tensor(  
249     preprocess_input_sequence(ex4_in,  
250         convert_str_to_tokens, SPECIAL_TOKENS)  
251 )  
252  
253 ex1_processed_expected = torch.tensor([14, 10, 0,  
254     3, 3, 3, 12, 10, 0, 6, 9, 6, 15])  
255 ex2_processed_expected = torch.tensor([14, 10, 0,  
256     6, 7, 3, 12, 10, 0, 6, 7, 5, 15])  
257 ex3_processed_expected = torch.tensor([14, 11, 0,  
258     2, 8, 6, 13, 11, 0, 0, 4, 4, 15])  
259 ex4_processed_expected = torch.tensor([14, 11, 0,  
260     4, 2, 0, 12, 10, 0, 3, 4, 2, 15])  
261  
262 ex1_out = torch.tensor(  
263     preprocess_input_sequence(ex1_out,  
264         convert_str_to_tokens, SPECIAL_TOKENS)  
265 )  
266 ex2_out = torch.tensor(  
267     preprocess_input_sequence(ex2_out,  
268         convert_str_to_tokens, SPECIAL_TOKENS)  
269 )  
270 ex3_out = torch.tensor(  
271     preprocess_input_sequence(ex3_out,  
272         convert_str_to_tokens, SPECIAL_TOKENS)  
273 )  
274 ex4_out = torch.tensor(  
275     preprocess_input_sequence(ex4_out,  
276         convert_str_to_tokens, SPECIAL_TOKENS)  
277 )  
278  
279 ex1_out_expected = torch.tensor([14, 10, 1, 0, 2,  
280     9, 15])  
281 ex2_out_expected = torch.tensor([14, 10, 1, 3, 4,  
282     8, 15])  
283 ex3_out_expected = torch.tensor([14, 11, 0, 2, 4,  
284     2, 15])
```

```
272 ex4_out_expected = torch.tensor([14, 11, 0, 0, 7,
273     8, 15])
274 print(
275     "preprocess input token error 1: ",
276     rel_error(ex1_processed_expected,
277     ex1_inp_preprocessed),
278 )
279 print(
280     "preprocess input token error 2: ",
281     rel_error(ex2_processed_expected,
282     ex2_inp_preprocessed),
283 )
284 print(
285     "preprocess input token error 3: ",
286     rel_error(ex3_processed_expected,
287     ex3_inp_preprocessed),
288 )
289 print(
290     "preprocess input token error 4: ",
291     rel_error(ex4_processed_expected,
292     ex4_inp_preprocessed),
293 )
294 print("\n")
295 print("preprocess output token error 1: ",
296     rel_error(ex1_out_expected, ex1_out))
297 print("preprocess output token error 2: ",
298     rel_error(ex2_out_expected, ex2_out))
299 print("preprocess output token error 3: ",
300     rel_error(ex3_out_expected, ex3_out))
301 print("preprocess output token error 4: ",
302     rel_error(ex4_out_expected, ex4_out))


---


295 #%% md
296 # Part II. Implementing Transformer building
297 blocks
298 Now that we have looked at the data, the task is
    to predict the output sequence (final result),
    something like `NEGATIVE 42` given the input
    sequence (of the arithmetic expression), something
    like `NEGATIVE 48 subtract NEGATIVE 6`.
```

299

300 In this section, we will look at implementing various building blocks used for implementing Transformer model. This will then be used to make Transformer encoder and decoder, which will ultimately lead us to implementing the complete Transformer model.

301 Each block will be implemented as a subclass of `nn.Module`; we will use PyTorch autograd to compute gradients, so we don't need to implement backward passes manually.

302

303 We will implement the following blocks, by referencing the original paper:

304

305 1. MultiHeadAttention Block

306 2. FeedForward Block

307 3. Layer Normalization

308 4. Positional Encoding block

309

310 We will then use these building blocks, combined with the input embedding layer to construct the Transformer Encoder and Decoder. We will start with MultiHeadAttention block, FeedForward Block, and Layer Normalization and look at Position encoding and input embedding later.

311

312 ****Note:**** One thing to keep in mind while implementing these blocks is that the shape of input and output Tensor from all these blocks we will be same. It always helps by checking the shapes of input and output tensors.

313 #%% md

314 MultiHeadAttention Block

315 #%% md

316 The MultiHead Attention block is a key component of the Transformer model. It allows the model to attend to different positions of the input sequence simultaneously.

317 #%% md

318 Transformers are sequence to sequence networks i.e

318 . , we get a sequence (for example a sentence in English) and get output a sequence (for example a sentence in Spanish). The input sequence are first transformed into embeddings as discussed in the RNN section and these embeddings are then passed through a Positional Encoding block. The resultant Embeddings are then transformed into three vectors, *query*, *key*, and *value* using learnable weights and we then use a Transformer Encoder and Decoder to get the final output sequence. For this section, we will assume that we have the *query*, *key*, and the *value* vector and work on them.

319

320 In the above figure, you can see that the Encoder has multihead attention block is right after these blocks. There is also a masked multihead attention in the deocoder but we will see that it's easy to implement the masked attention when we have implemented the basic MultiHeadAttention block.

321 To implement the basic MultiheadAttention block, we will first implement the Self Attention block and see that MultiHeadAttention can be implemented as a direct extension of the Self Attention block .

322

323 ## Self Attention Block

324

325 Taking inspiration from information retrieval paradigm, Transformers have this notion of *query*, *key*, and *value* where given a *query* we try extract information from *key*-*value* pairs. Moving along those lines, we perform this mathematically by taking the weighted sum of *values* for each *query*, where weight is computed by dot product of *query* and the *key*. More precisely, for each query we compute the dot product with all the keys and then use the scalar output of those dot products as weights to find the weighted sum of *values*. Note that before

```

325 finding the weighted sum, we also apply softmax
      function to the weights vector. Lets start with
      implementing of Attention Block that takes input
      as *query*, *key*, and *value* vectors and returns
      a Tensor, that is weighted sum of the *values*.
326
327 For this section, you need to implement three
      functions, `scaled_dot_product_two_loop_single`, `scaled_dot_product_two_loop_batch` and `scaled_dot_product_no_loop_batch` inside the
      transformers.py file. This might look very similar
      to the `dot_product_attention` in the RNN
      notebook but there is a subtle difference in the
      inputs. You should see the errors of the order
      less than 1e-5


---


328 #%%
329 from transformers import (
330     scaled_dot_product_two_loop_single,
331     scaled_dot_product_two_loop_batch,
332     scaled_dot_product_no_loop_batch,
333 )


---


334 #%%
335 N = 2 # Number of sentences
336 K = 5 # Number of words in a sentence
337 M = 4 # feature dimension of each word embedding
338
339 query = torch.linspace(-0.4, 0.6, steps=K * M).
    reshape(K, M) # **to_double_cuda
340 key = torch.linspace(-0.8, 0.5, steps=K * M).
    reshape(K, M) # **to_double_cuda
341 value = torch.linspace(-0.3, 0.8, steps=K * M).
    reshape(K, M) # **to_double_cuda
342
343 y = scaled_dot_product_two_loop_single(query, key
    , value)
344 y_expected = torch.tensor(
345     [
346         [0.08283, 0.14073, 0.19862, 0.25652],
347         [0.13518, 0.19308, 0.25097, 0.30887],
348         [0.18848, 0.24637, 0.30427, 0.36216],
349         [0.24091, 0.29881, 0.35670, 0.41460],

```

```

350 [0.29081, 0.34871, 0.40660, 0.46450],
351 ]
352 ).to(torch.float32)
353 print("scaled_dot_product_two_loop_single error: "
, rel_error(y_expected, y))
354 #%%
355 N = 2 # Number of sentences
356 K = 5 # Number of words in a sentence
357 M = 4 # feature dimension of each word embedding
358
359 query = torch.linspace(-0.4, 0.6, steps=N * K * M
).reshape(N, K, M) # **to_double_cuda
360 key = torch.linspace(-0.8, 0.5, steps=N * K * M).
reshape(N, K, M) # **to_double_cuda
361 value = torch.linspace(-0.3, 0.8, steps=N * K * M
).reshape(N, K, M) # *to_double_cuda
362
363 y = scaled_dot_product_two_loop_batch(query, key,
value)
364 y_expected = torch.tensor(
365 [
366 [
367 [-0.09603, -0.06782, -0.03962, -0.
01141],
368 [-0.08991, -0.06170, -0.03350, -0.
00529],
369 [-0.08376, -0.05556, -0.02735, 0.00085
],
370 [-0.07760, -0.04939, -0.02119, 0.00702
],
371 [-0.07143, -0.04322, -0.01502, 0.01319
],
372 ],
373 [
374 [0.49884, 0.52705, 0.55525, 0.58346],
375 [0.50499, 0.53319, 0.56140, 0.58960],
376 [0.51111, 0.53931, 0.56752, 0.59572],
377 [0.51718, 0.54539, 0.57359, 0.60180],
378 [0.52321, 0.55141, 0.57962, 0.60782],
379 ],
380 ]

```

```
381 ).to(torch.float32)
382 print("scaled_dot_product_two_loop_batch error: "
383     , rel_error(y_expected, y))
384 #%%
385 N = 2 # Number of sentences
386 K = 5 # Number of words in a sentence
387 M = 4 # feature dimension of each word embedding
388
389 query = torch.linspace(-0.4, 0.6, steps=N * K * M
390     ).reshape(N, K, M) # **to_double_cuda
391 key = torch.linspace(-0.8, 0.5, steps=N * K * M).
392     reshape(N, K, M) # **to_double_cuda
393 value = torch.linspace(-0.3, 0.8, steps=N * K * M
394     ).reshape(N, K, M) # *to_double_cuda
395
396 y, _ = scaled_dot_product_no_loop_batch(query, key,
397     , value)
398
399 y_expected = torch.tensor(
400     [
401         [
402             [
403                 [
404                     [
405                         [-0.09603, -0.06782, -0.03962, -0.
406                          01141],
407                         [-0.08991, -0.06170, -0.03350, -0.
408                          00529],
409                         [-0.08376, -0.05556, -0.02735, 0.00085
410                         ],
411                         [-0.07760, -0.04939, -0.02119, 0.00702
412                         ],
413                         [-0.07143, -0.04322, -0.01502, 0.01319
414                         ],
415                         ],
416                         [
417                             [0.49884, 0.52705, 0.55525, 0.58346],
418                             [0.50499, 0.53319, 0.56140, 0.58960],
419                             [0.51111, 0.53931, 0.56752, 0.59572],
420                             [0.51718, 0.54539, 0.57359, 0.60180],
421                             [0.52321, 0.55141, 0.57962, 0.60782],
422                             ],
423                         ]
424                     ]
425                 ]
426             ]
427         ]
428     ]
```

```

412 ).to(torch.float32)
413
414 print("scaled_dot_product_no_loop_batch error: ",
    rel_error(y_expected, y))


---


415 #%% md
416 ## Observing time complexity:
417
418 As Transformers are infamous for their time
    complexity that depends on the size of the input
    sequence.
419 We can verify this now that we have implemented
    `self_attention_no_loop`.
420 Run the cells below: the first has a sequence
    length of 256 and the second one has a sequence
    length of 512. You should roughly be 4 times
    slower with sequence length 512, hence showing
    that compleixity of the transformers increase
    quadratically with respect to increase in the in
    sequence length.
421 The `%timeit` lines may take several seconds to
    run.


---


422 #%%
423 N = 64
424 K = 256 # defines the input sequence length
425 M = emb_size = 2048
426 dim_q = dim_k = 2048
427 query = torch.linspace(-0.4, 0.6, steps=N * K * M
    ).reshape(N, K, M) # **to_double_cuda
428 key = torch.linspace(-0.8, 0.5, steps=N * K * M).
    reshape(N, K, M) # **to_double_cuda
429 value = torch.linspace(-0.3, 0.8, steps=N * K * M
    ).reshape(N, K, M) # *to_double_cuda
430
431 %timeit -n 5 -r 2 y =
    scaled_dot_product_no_loop_batch(query, key, value
    )


---


432 #%%
433 N = 64
434 K = 512 # defines the input requence length
435 M = emb_size = 2048
436 dim_q = dim_k = 2048

```

```

437 query = torch.linspace(-0.4, 0.6, steps=N * K * M
    ).reshape(N, K, M) # **to_double_cuda
438 key = torch.linspace(-0.8, 0.5, steps=N * K * M).
    reshape(N, K, M) # **to_double_cuda
439 value = torch.linspace(-0.3, 0.8, steps=N * K * M
    ).reshape(N, K, M) # **to_double_cuda
440
441 %timeit -n 5 -r 2 y =
    scaled_dot_product_no_loop_batch(query, key, value
)
442 #%% md
443 Now that we have implemented `scaled_dot_product_no_loop_batch`, lets implement `SingleHeadAttention`, that will serve as a building block for the `MultiHeadAttention` block. For this exercise, we have made a `SingleHeadAttention` class that inherits from `nn.module` class of Pytorch. You need to implement the `__init__` and the `forward` functions inside `Transformers.py`
444 #%% md
445 Run the following cells to test your implementation of `SelfAttention` layer. We have also written code to check the backward pass using pytorch autograd API in the following cell. You should expect the error to be less than 1e-5
446 #%%
447 from transformers import SelfAttention
448 #%%
449 reset_seed(0)
450 N = 2
451 K = 4
452 M = emb_size = 4
453 dim_q = dim_k = 4
454 atten_single = SelfAttention(emb_size, dim_q,
    dim_k)
455
456 for k, v in atten_single.named_parameters():
457     # print(k, v.shape) # uncomment this to see
        the weight shape
458     v.data.copy_(torch.linspace(-1.4, 1.3, steps=v

```

```
458 .numel()).reshape(*v.shape))
459
460 query = torch.linspace(-0.4, 0.6, steps=N * K * M
, requires_grad=True).reshape(
461     N, K, M
462 ) # **to_double_cuda
463 key = torch.linspace(-0.8, 0.5, steps=N * K * M,
requires_grad=True).reshape(
464     N, K, M
465 ) # **to_double_cuda
466 value = torch.linspace(-0.3, 0.8, steps=N * K * M
, requires_grad=True).reshape(
467     N, K, M
468 ) # **to_double_cuda
469
470 query,retain_grad()
471 key,retain_grad()
472 value,retain_grad()
473
474 y_expected = torch.tensor(
475     [
476         [
477             [-1.10382, -0.37219, 0.35944, 1.09108
478         ],
479             [-1.45792, -0.50067, 0.45658, 1.41384
480         ],
481             [-1.74349, -0.60428, 0.53493, 1.67414
482         ],
483             [-1.92584, -0.67044, 0.58495, 1.84035
484         ],
485         [
486             [-4.59671, -1.63952, 1.31767, 4.27486
487         ],
```

```
488     ]
489 )
490
491 dy_expected = torch.tensor(
492     [
493         [
494             [-0.09084, -0.08961, -0.08838, -0.
495              08715],
496              [0.69305, 0.68366, 0.67426, 0.66487],
497              [-0.88989, -0.87783, -0.86576, -0.
498              85370],
499              [0.25859, 0.25509, 0.25158, 0.24808],
500          ],
501          [
502              [-0.05360, -0.05287, -0.05214, -0.
503              05142],
504              [0.11627, 0.11470, 0.11312, 0.11154],
505              [-0.01048, -0.01034, -0.01019, -0.
506              01005],
507              [-0.03908, -0.03855, -0.03802, -0.
508              03749],
509          ],
510      ]
511 )
512
513
514 print("SelfAttention error: ", rel_error(
515     y_expected, y))
516 print("SelfAttention error: ", rel_error(
517     dy_expected, query_grad))
518
519 #%% md
520
521 We have implemented the `SingleHeadAttention`  
block which brings us very close to implementing  
`MultiHeadAttention`. We will now see that this  
can be achieved by manipulating the shapes of  
input tensors based on number of heads in the
```

```

517 Multi-Attention block. We design a network that
    uses multiple SingleHeadAttention blocks on the
    same input to compute the output tensors and
    finally concatenate them to generate a single
    output. This is not the implementation used in
    practice as it forces you to initialize multiple
    layers but we use it here for simplicity.
    Implement MultiHeadAttention block in the
    transformers.py` file by using the
    SingleHeadAttention block.
518 #%% md
519 Run the following cells to test your
    MultiHeadAttention` layer. Again, as
    SelfAttention`, we have used pytorch autograd API
    to test the backward pass. You should expect error
    values below 1e-5.
520 #%%
521 from transformers import MultiHeadAttention
522 #%%
523 reset_seed(0)
524 N = 2
525 num_heads = 2
526 K = 4
527 M = inp_emb_size = 4
528 out_emb_size = 8
529 atten_multihead = MultiHeadAttention(num_heads,
    inp_emb_size, out_emb_size)
530
531 for k, v in atten_multihead.named_parameters():
532     # print(k, v.shape) # uncomment this to see
    the weight shape
533     v.data.copy_(torch.linspace(-1.4, 1.3, steps=v
        .numel()).reshape(*v.shape))
534
535 query = torch.linspace(-0.4, 0.6, steps=N * K * M
    , requires_grad=True).reshape(
536     N, K, M
537 ) # **to_double_cuda
538 key = torch.linspace(-0.8, 0.5, steps=N * K * M,
    requires_grad=True).reshape(
539     N, K, M

```

```

540 ) # **to_double_cuda
541 value = torch.linspace(-0.3, 0.8, steps=N * K * M
, requires_grad=True).reshape(
542     N, K, M
543 ) # *to_double_cuda
544
545 query,retain_grad()
546 key,retain_grad()
547 value,retain_grad()
548
549 y_expected = torch.tensor(
550     [
551         [
552             [-0.23104, 0.50132, 1.23367, 1.96603],
553             [0.68324, 1.17869, 1.67413, 2.16958],
554             [1.40236, 1.71147, 2.02058, 2.32969],
555             [1.77330, 1.98629, 2.19928, 2.41227],
556         ],
557         [
558             [6.74946, 5.67302, 4.59659, 3.52015],
559             [6.82813, 5.73131, 4.63449, 3.53767],
560             [6.86686, 5.76001, 4.65315, 3.54630],
561             [6.88665, 5.77466, 4.66268, 3.55070],
562         ],
563     ]
564 )
565 dy_expected = torch.tensor(
566     [[[ 0.56268, 0.55889, 0.55510, 0.55131],
567       [ 0.43286, 0.42994, 0.42702, 0.42411],
568       [ 2.29865, 2.28316, 2.26767, 2.25218],
569       [ 0.49172, 0.48841, 0.48509, 0.48178
570     ],
571     [[ 0.25083, 0.24914, 0.24745, 0.24576],
572       [ 0.14949, 0.14849, 0.14748, 0.14647],
573       [-0.03105, -0.03084, -0.03063, -0.03043],
574       [-0.02082, -0.02068, -0.02054, -0.02040
575     ]]
576 )
577 y = atten_multihead(query, key, value)

```

```

578 dy = torch.randn(*y.shape) # , **to_double_cuda
579
580 y.backward(dy)
581 query_grad = query.grad
582 print("MultiHeadAttention error: ", rel_error(
    y_expected, y))
583 print("MultiHeadAttention error: ", rel_error(
    dy_expected, query_grad))


---


584 #%% md
585 #### LayerNormalization
586 #%% md
587 Layer Normalization is used throughout the
Transformer model. It normalizes the inputs across
the features for each sample independently,
helping to stabilize training.


---


588 #%% md
589 We implemented BatchNorm while working with CNNs.
One of the problems of BatchNorm is its dependency
on the complete batch which might not give
good results when the batch size is small. Ba et
al proposed `LayerNormalization` that takes into
account these problems and has become a standard
in sequence-to-sequence tasks. In this section, we
will implement `LayerNormalization`. Another nice
quality of `LayerNormalization` is that as it
depends on individual time steps or each element
of the sequence, it can be parallelized and the
test time runs in a similar manner hence making it
better implementation wise. Again, you have to
only implement the forward pass and the backward
pass will be taken care by Pytorch autograd.
Implement the `LayerNormalization` class in
`transformers.py`, you should expect the error
below 1e-5


---


590 #%%
591 from transformers import LayerNormalization


---


592 #%%
593 reset_seed(0)
594 N = 2
595 K = 4
596 norm = LayerNormalization(K)

```

```

597 inp = torch.linspace(-0.4, 0.6, steps=N * K,
   requires_grad=True).reshape(N, K)
598
599 inp.retain_grad()
600 y = norm(inp)
601
602 y_expected = torch.tensor(
603     [[-1.34164, -0.44721, 0.44721, 1.34164], [-1.
604      34164, -0.44721, 0.44721, 1.34164]]
604 )
605
606 dy_expected = torch.tensor(
607     [[ 5.70524, -2.77289, -11.56993,  8.63758],
608      [ 2.26242, -4.44330,  2.09933,  0.
609      08154]]
609 )
610
611 dy = torch.randn(*y.shape)
612 y.backward(dy)
613 inp_grad = inp.grad
614
615 print("LayerNormalization error: ", rel_error(
   y_expected, y))
616 print("LayerNormalization grad error: ", rel_error(
   dy_expected, inp_grad))


---


617 #%% md
618 ### FeedForward Block


---


619 #%% md
620 The FeedForward Block is used in both the Encoder  

   and Decoder of the Transformer. It consists of two  

   linear transformations with a ReLU activation in  

   between.


---


621 #%% md
622 Next, we will implement the `Feedforward` block.  

   These are used in both the Encoder and Decoder  

   network of the Transformer and they consist of  

   stacked MLP and ReLU layers. In the overall  

   architecture, the output of `MultiHeadAttention`  

   is fed into the `FeedForward` block. Implement the  

   `FeedForwardBlock` inside `transformers.py` and  

   execute the following cells to check your

```

```
622 implementation. You should expect the errors below
623      1e-5
624 #%%
624 from transformers import FeedForwardBlock
625 #%%
626 reset_seed(0)
627 N = 2
628 K = 4
629 M = emb_size = 4
630
631 ff_block = FeedForwardBlock(emb_size, 2 * emb_size
632 )
633 for k, v in ff_block.named_parameters():
634     v.data.copy_(torch.linspace(-1.4, 1.3, steps=v
635 .numel()).reshape(*v.shape))
636 inp = torch.linspace(-0.4, 0.6, steps=N * K,
637 requires_grad=True).reshape(
638     N, K
639 )
639 inp,retain_grad()
640 y = ff_block(inp)
641
642 y_expected = torch.tensor(
643     [[-2.46161, -0.71662, 1.02838, 2.77337], [-7.
644 56084, -1.69557, 4.16970, 10.03497]]
645 )
646 dy_expected = torch.tensor(
647     [[0.55105, 0.68884, 0.82662, 0.96441], [0.
648 30734, 0.31821, 0.32908, 0.33996]]
649 )
650 dy = torch.randn(*y.shape)
651 y.backward(dy)
652 inp_grad = inp.grad
653
654 print("FeedForwardBlock error: ", rel_error(
655     y_expected, y))
655 print("FeedForwardBlock error: ", rel_error(
```

```
655 dy_expected, inp_grad))
656 #%% md
657 Now, if you look back to the original paper,
Attention is all you Need, then, we are almost
done with the building blocks of a transformer.
What's left is:
658
659 - Encapsulating the building blocks into Encoder
Block
660 - Encapsulating the building blocks into Decoder
Block
661 - Handling the input data preprocessing and
positional encoding.
662
663 We will first look at implementing the Encoder
Block and Decoder block. The positional encoding
is a non learnable embedding and we can treat it
as a preprocessing step in our DataLoader.
664 #%% md
665 The Encoder Block combines MultiHead Attention,
Layer Normalization, and FeedForward layers with
residual connections. Each encoder block processes
the input sequence and passes it to the next
layer.
666 #%% md
667 As shown in the figure above, the encoder block
takes it inputs three tensors. We will assume that
we have those three tensors, query, key, and
value. Run the cell below to check your
implementation of the EncoderBlock. You should
expect the errors below 1e-5
668 #%%
669 from transformers import EncoderBlock
670 #%%
671 reset_seed(0)
672 N = 2
673 num_heads = 2
674 emb_dim = K = 4
675 feedforward_dim = 8
676 M = inp_emb_size = 4
677 out_emb_size = 8
```

```
678 dropout = 0.2
679
680 enc_seq_inp = torch.linspace(-0.4, 0.6, steps=N *
681     K * M, requires_grad=True).reshape(
682     N, K, M
683 ) # **to_double_cuda
684
685 enc_block = EncoderBlock(num_heads, emb_dim,
686     feedforward_dim, dropout)
687
688 for k, v in enc_block.named_parameters():
689     # print(k, v.shape) # uncomment this to see
690     # the weight shape
691     v.data.copy_(torch.linspace(-1.4, 1.3, steps=v
692         .numel()).reshape(*v.shape))
693
694 encoder_out1_expected = torch.tensor(
695     [[[ 0.00000, -0.31357,  0.69126,  0.00000],
696       [ 0.42630, -0.25859,  0.72412,  3.87013],
697       [ 0.00000, -0.31357,  0.69126,  3.89884],
698       [ 0.47986, -0.30568,  0.69082,  3.90563
699     ]],
700
701     [[[ 0.00000, -0.31641,  0.69000,  3.89921],
702       [ 0.47986, -0.30568,  0.69082,  3.90563],
703       [ 0.47986, -0.30568,  0.69082,  3.90563],
704       [ 0.51781, -0.30853,  0.71598,  3.85171
705     ]])
706
707 encoder_out1 = enc_block(enc_seq_inp)
708 print("EncoderBlock error 1: ", rel_error(
709     encoder_out1, encoder_out1_expected))
710
711 N = 2
712 num_heads = 1
713 emb_dim = K = 4
714 feedforward_dim = 8
715 M = inp_emb_size = 4
716 out_emb_size = 8
717 dropout = 0.2
```

```

712
713 enc_seq_inp = torch.linspace(-0.4, 0.6, steps=N *
    K * M, requires_grad=True).reshape(
714     N, K, M
715 ) # **to_double_cuda
716
717 enc_block = EncoderBlock(num_heads, emb_dim,
    feedforward_dim, dropout)
718
719 for k, v in enc_block.named_parameters():
720     # print(k, v.shape) # uncomment this to see
    the weight shape
721     v.data.copy_(torch.linspace(-1.4, 1.3, steps=v
        .numel()).reshape(*v.shape))
722
723 encoder_out2_expected = torch.tensor(
    [[[ 0.42630, -0.00000,  0.72412,  3.87013],
724      [ 0.49614, -0.31357,  0.00000,  3.89884],
725      [ 0.47986, -0.30568,  0.69082,  0.00000],
726      [ 0.51654, -0.32455,  0.69035,  3.89216
    ]],
728
729      [[ 0.47986, -0.30568,  0.69082,  0.00000],
730      [ 0.49614, -0.31357,  0.69126,  3.89884],
731      [ 0.00000, -0.30354,  0.76272,  3.75311],
732      [ 0.49614, -0.31357,  0.69126,  3.89884
    ]]]
733 )
734 encoder_out2 = enc_block(enc_seq_inp)
735 print("EncoderBlock error 2: ", rel_error(
    encoder_out2, encoder_out2_expected))


---


736 %% md
737 Great! You're almost done with the implementation
    of the Transformer model.


---


738 %% md
739 ### Decoder Block
740
741 The Decoder Block is similar to the Encoder but
    includes an additional masked self-attention layer
    and a cross-attention layer that attends to the
    encoder output.

```

742 #%% md

743 Now, we will look at the implementation of the decoder. In the class we learned about encoder only model that can be used for tasks like sequence classification but for more complicated tasks like sequence to sequence we need a decoder network that can transform the output of the encoder to a target sequence. This kind of architecture is important in tasks like language translation where we have a sequence as input and a sequence as output. This decoder takes the input from the encoder and the previous generated value to generate the next value. During training, we use a Mask on the input so that the decoder network can't look ahead in the future and during inference we sequentially process the data.

744

745 Before moving to implementing the Decoder Block, we should pay attention to the figure above. It says a "Masked MultiHead Attention" which actually prevents the decoder from looking ahead into the future. Lets understand with an example here. We have an expression as `BOS POSITIVE 01 add POSITIVE 00 EOS`, i.e. `1+0` that gives output as `BOS POSITIVE 01 EOS`, i.e. `+1`. Lets focus on the output sequence here. This is a sequence of length 5 (after applying our preprocessing code) and will get transformed into *key*, *query*, and *value* matrix of dimension \$5\times128\$, \$5\times128\$ and \$5\times128\$ respectively, where 128 is the embedding dimension of the Transformer. Now, while training, we input these vectors in the `self_attention_no_loop_batch` without mask. It will compute the dot product between *query* and *key* to generate a \$5\times5\$ matrix where the first row (shape \$1\times5\$) of that matrix tells us how much the word `EOS` is related with `EOS`, `POSITIVE`, `0`, `1`, and `EOS`. This means that it will use the weights of all these tokens to learn the final sequence that is to be predicted. This is okay when we are training the

745 model but what happens when we perform inference?
 We start with a brand new expression, input this
 expression in the encoder but this time we only
 have the first starting token `EOS` for decoder
 and we don't know about the rest of the tokens in
 the sequence. Hence, a solution to this problem is
 to mask the weights inside the function
 `self_attention_no_loop_batch` for only the decoder
 part. This masking should prevent the decoder
 from accessing the future or next elements.

746

747 We will now look at how to generate this mask for
 a given sequence. Then, you should also update the
 `self_attention_no_loop_batch` to use the mask
 variable appropriately. Implement the
 `get_subsequent_mask`,
 `self_attention_no_loop_batch` with mask inside
 `transformers.py` file

748 #%%

```
749 from transformers import get_subsequent_mask
```

750

```
751 reset_seed(0)
```

```
752 seq_len_enc = K = 4
```

```
753 M = inp_emb_size = 3
```

754

```
755 inp_sequence = torch.linspace(-0.4, 0.6, steps=K  

    * M, requires_grad=True).reshape(
```

```
756     K, M
```

```
757 ) # **to_double_cuda
```

758

```
759 mask_expected = torch.tensor(
```

```
760     [
```

```
761         [[False, True, True], [False, False, True  

    ], [False, False, False]],
```

```
762         [[False, True, True], [False, False, True  

    ], [False, False, False]],
```

```
763         [[False, True, True], [False, False, True  

    ], [False, False, False]],
```

```
764         [[False, True, True], [False, False, True  

    ], [False, False, False]],
```

```
765     ]
```

```
766 )
767 mask_predicted = get_subsequent_mask(inp_sequence)
768 print(
769     "get_subsequent_mask error: ", rel_error(
770         mask_predicted.int(), mask_expected.int()))
770 )
771 #%%
772 from transformers import
    scaled_dot_product_no_loop_batch
773 #%%
774 reset_seed(0)
775 N = 4
776 K = 3
777 M = 3
778
779 query = torch.linspace(-0.4, 0.6, steps=K * M * N
    , requires_grad=True).reshape(N, K, M)
780 key = torch.linspace(-0.1, 0.2, steps=K * M * N,
    requires_grad=True).reshape(N, K, M)
781 value = torch.linspace(0.4, 0.8, steps=K * M * N,
    requires_grad=True).reshape(N, K, M)
782
783 y_expected = torch.tensor(
784     [
785         [
786             [0.40000, 0.41143, 0.42286],
787             [0.41703, 0.42846, 0.43989],
788             [0.43408, 0.44551, 0.45694],
789         ],
790         [
791             [0.50286, 0.51429, 0.52571],
792             [0.51999, 0.53142, 0.54285],
793             [0.53720, 0.54863, 0.56006],
794         ],
795         [
796             [0.60571, 0.61714, 0.62857],
797             [0.62294, 0.63437, 0.64580],
798             [0.64032, 0.65175, 0.66318],
799         ],
800         [
801             [0.70857, 0.72000, 0.73143],
```

```
802 [0.72590, 0.73733, 0.74876],  
803 [0.74344, 0.75487, 0.76630],  
804 ],  
805 ]  
806 )  
807 y_predicted, _ = scaled_dot_product_no_loop_batch(  
    query, key, value, mask_expected)  
808  
809 print("scaled_dot_product_no_loop_batch error: ",  
    rel_error(y_expected, y_predicted))  
810 #%% md  
811 Lets finally implement the decoder block now that  
we have all the required tools to implement it.  
Fill in the init function and the forward pass of  
the `DecoderBlock` inside `transformers.py`. Run  
the following cells to check your implementation  
of the `DecoderBlock`. You should expect the  
errors below 1e-5.  
812 #%%  
813 from transformers import DecoderBlock  
814 #%%  
815 reset_seed(0)  
816 N = 2  
817 num_heads = 2  
818 seq_len_enc = K1 = 4  
819 seq_len_dec = K2 = 2  
820 feedforward_dim = 8  
821 M = emb_dim = 4  
822 out_emb_size = 8  
823 dropout = 0.2  
824  
825 dec_inp = torch.linspace(-0.4, 0.6, steps=N * K1  
    * M, requires_grad=True).reshape(  
826     N, K1, M  
827 )  
828 enc_out = torch.linspace(-0.4, 0.6, steps=N * K2  
    * M, requires_grad=True).reshape(  
829     N, K2, M  
830 )  
831 dec_block = DecoderBlock(num_heads, emb_dim,  
    feedforward_dim, dropout)
```

```

832
833 for k, v in dec_block.named_parameters():
834     # print(k, v.shape) # uncomment this to see
835     # the weight shape
836     v.data.copy_(torch.linspace(-1.4, 1.3, steps=v
837     .numel()).reshape(*v.shape))
838
839 dec_out_expected = torch.tensor(
840     [[[ 0.50623, -0.32496,  0.00000,  0.00000],
841       [ 0.00000, -0.31690,  0.76956,  3.72647],
842       [ 0.49014, -0.32809,  0.66595,  3.93773],
843       [ 0.00000, -0.00000,  0.68203,  3.90856
844     ],
845     [[ 0.51042, -0.32787,  0.68093,  3.90848],
846     [ 0.00000, -0.31637,  0.72275,  3.83122],
847     [ 0.64868, -0.00000,  0.77715,  0.00000],
848     [ 0.00000, -0.33105,  0.66565,  3.93602
849   ]]
850 )
851 dec_out1 = dec_block(dec_inp, enc_out)
852 print("DecoderBlock error: ", rel_error(dec_out1,
853                                         dec_out_expected))
854
855 N = 2
856 num_heads = 2
857 seq_len_enc = K1 = 4
858 seq_len_dec = K2 = 4
859 feedforward_dim = 4
860 emb_dim = 4
861 out_emb_size = 8
862 dropout = 0.2
863
864 dec_inp = torch.linspace(-0.4, 0.6, steps=N * K1
865     * M, requires_grad=True).reshape(
866     N, K1, M
867 )
868 enc_out = torch.linspace(-0.4, 0.6, steps=N * K2
869     * M, requires_grad=True).reshape(
870     N, K2, M

```

```

866 )
867 dec_block = DecoderBlock(num_heads, emb_dim,
   feedforward_dim, dropout)
868
869 for k, v in dec_block.named_parameters():
870     # print(k, v.shape) # uncomment this to see
     the weight shape
871     v.data.copy_(torch.linspace(-1.4, 1.3, steps=v
       .numel()).reshape(*v.shape))
872
873
874 dec_out_expected = torch.tensor(
875     [[[ 0.46707, -0.31916,  0.66218,  3.95182],
876       [ 0.00000, -0.31116,  0.66325,  0.00000],
877       [ 0.44538, -0.32419,  0.64068,  3.98847],
878       [ 0.49012, -0.31276,  0.68795,  3.90610
879     ],
880     [[ 0.45800, -0.33023,  0.64106,  3.98324],
881     [ 0.45829, -0.31487,  0.66203,  3.95529],
882     [ 0.59787, -0.00000,  0.72361,  0.00000],
883     [ 0.70958, -0.37051,  0.78886,  3.63179
884   ]]
885 )
886 dec_out2 = dec_block(dec_inp, enc_out)
887 print("DecoderBlock error: ", rel_error(dec_out2,
     dec_out_expected))
888 #%% md
889 Based on the implementation of `EncoderBlock` and
   `DecoderBlock`, we have implemented the `Encoder`
   and `Decoder` networks for you in transformers.
   py. You should be able to understand the input and
   outputs of these Encoder and Decoder blocks.
   Implement the Transformer block inside transformer
   .py using these networks.
890 #%% md
891 ## Part III: Data loader
892 In this part, we will have a look at creating the
   final data loader for the task, that can be used
   to train the Transformer model. This will comprise

```

```
892 of two things:  
893  
894 - Implement Positional Encoding  
895 - Create a dataloader using the  
     preprocess_input_sequence` function that we created  
     in Part I.  
896 #%% md  
897 Lets start with implementing the Positional  
     Encoding for the input. The positional encodings  
     make the Transformers positionally aware about  
     sequences. These are usually added to the input  
     and hence should be same shape as input. As these  
     are not learnable, they remain constant  
     throughout the training process. For this reason  
     , we can look at it as a pre-processing step that'  
     s done on the input. Our strategy here would be to  
     implement positional encoding function and use it  
     later while creating DataLoader for the toy  
     dataset.  
898  
899 Lets look at the simplest kind of positional  
     encoding, i.e. for a sequence of length K, assign  
     the nth element in the sequence a value of n/K,  
     where n starts from 0. Implement the  
     position_encoding_simple inside `transformers.py  
     `. You should expect error less than 1e-9 here.  
900 #%% md  
901 ### Simple positional encoding  
902 #%%  
903 from transformers import position_encoding_simple  
904  
905 reset_seed(0)  
906 K = 4  
907 M = emb_size = 4  
908  
909 y = position_encoding_simple(K, M)  
910 y_expected = torch.tensor(  
911     [  
912         [  
913             [0.00000, 0.00000, 0.00000, 0.00000],  
914             [0.25000, 0.25000, 0.25000, 0.25000],
```

```

915 [0.50000, 0.50000, 0.50000, 0.50000],
916 [0.75000, 0.75000, 0.75000, 0.75000],
917 ]
918 ]
919 )
920
921 print("position_encoding_simple error: ",
922     rel_error(y, y_expected))
923 K = 5
924 M = emb_size = 3
925
926
927 y = position_encoding_simple(K, M)
928 y_expected = torch.tensor(
929     [
930         [
931             [0.00000, 0.00000, 0.00000],
932             [0.20000, 0.20000, 0.20000],
933             [0.40000, 0.40000, 0.40000],
934             [0.60000, 0.60000, 0.60000],
935             [0.80000, 0.80000, 0.80000],
936         ]
937     ]
938 )
939 print("position_encoding_simple error: ",
940     rel_error(y, y_expected))
941 ### Sinusoid positional encoding
942 #% md
943 Now that we have looked at a simple positional
encoding, we can see one major drawback, which is
that if the sequence length gets larger, the
difference between two consecutive positional
encodings becomes smaller and smaller and it in
turn defeats a purpose of positional awareness, as
there is very small difference in two consecutive
positions. Another issue is that for each position
we replicated it along embedding dimension, hence
introducing redundancy which might not help the
network in learning anything new. There could be

```

```

943 different tricks that can be used to make a
    positional encoding that could solve these
    problems.
944
945 Lets look at more mature version of a positonal
    encoding that uses a combination of sines and
    cosines function, also called sinusoid. This is
    also the positional encoding used in the original
    Transformer paper. For each element in the
    sequence (length K) with position $p$ and
    embedding (dimension M) positon $i$, we can define
    the positional encoding as:
946
947 
$$\text{PE}_{(p, 2i)} = \sin\left(\frac{p}{10000^a}\right)$$

948 
$$\text{PE}_{(p, 2i+1)} = \cos\left(\frac{p}{10000^a}\right)$$

949
950 
$$\text{Where } a = \lfloor \frac{2i}{M} \rfloor \text{ and } M \text{ is the Embedding dimension of the Transformer}$$

951
952 Here, $p$ remains constant for a position in the
    sequence and we assign alternating sines and
    cosines along the embedding dimension.
953
954 Implement the fucntion `position_encoding` inside
    `transformers.py`. You should expect errors below
     $1e-5$ .
955 #%%
956 from transformers import
    position_encoding_sinusoid
957
958 reset_seed(0)
959 K = 4
960 M = emb_size = 4
961
962 y1 = position_encoding_sinusoid(K, M)
963 y_expected = torch.tensor(
964     [
965         [

```

```
966 [0.00000, 1.00000, 0.00000, 1.00000],  
967 [0.84147, 0.54030, 0.84147, 0.54030],  
968 [0.90930, -0.41615, 0.90930, -0.41615  
],  
969 [0.14112, -0.98999, 0.14112, -0.98999  
],  
970 ]  
971 ]  
972 )  
973  
974 print("position_encoding error: ", rel_error(y1,  
y_expected))  
975  
976 K = 5  
977 M = emb_size = 3  
978  
979  
980 y2 = position_encoding_sinusoid(K, M)  
981 y_expected = torch.tensor(  
982 [  
983 [  
984 [0.00000, 1.00000, 0.00000],  
985 [0.84147, 0.54030, 0.84147],  
986 [0.90930, -0.41615, 0.90930],  
987 [0.14112, -0.98999, 0.14112],  
988 [-0.75680, -0.65364, -0.75680],  
989 ]  
990 ]  
991 )  
992 print("position_encoding error: ", rel_error(y2,  
y_expected))  
993 #%% md  
994 ### Constructing the DataLoader for the toy  
dataset  
995 #%% md  
996 Now we will use the implemented positonal  
encodings to construct a DataLoader in Pytorch.  
The function of a data loader is to return a  
batch for training/validation. We first make a  
Dataset class that gives us a single element in  
the batch and then use a DataLoader to wrap the
```

```
996 dataset. We inherit the Dataset from `torch.utils  
997 .data.Dataset` class. This class consists of two  
998 important functions that you'd change depending  
999 on your usecase (for e.g. the upcoming project  
1000 !). The first function is `__init__`, this  
1001 consists of the components that are *static*, in  
1002 other words, these are the variables that won't  
1003 change when we want the next element from the  
1004 complete data. The second fucntion is `  
1005 __getitem__` which contains the core  
1006 functionality of the final dataloader.  
1007  
1008 To get the final dataloader, we wrap the `  
1009 train_data` and `test_data` in `torch.utils.data.  
1010 DataLoader` class  
1011 #%%  
1012 from sklearn.model_selection import  
1013     train_test_split  
1014 from transformers import AddSubDataset  
1015  
1016 BATCH_SIZE = 16  
1017  
1018 X, y = data["inp_expression"], data["  
1019     out_expression"]  
1020  
1021 X_train, X_test, y_train, y_test =  
1022     train_test_split(X, y, test_size=0.1,  
1023     random_state=0)  
1024  
1025 train_data = AddSubDataset(  
1026     X_train,  
1027     y_train,  
1028     convert_str_to_tokens,  
1029     SPECIAL_TOKENS,  
1030     32,  
1031     position_encoding_simple,  
1032 )  
1033 valid_data = AddSubDataset(  
1034     X_test, y_test, convert_str_to_tokens,  
1035     SPECIAL_TOKENS, 32, position_encoding_simple  
1036 )
```

```
1020
1021 train_loader = torch.utils.data.DataLoader(
1022     train_data, batch_size=BATCH_SIZE, shuffle=
1023     False, drop_last=True
1024 )
1025 valid_loader = torch.utils.data.DataLoader(
1026     valid_data, batch_size=BATCH_SIZE, shuffle=
1027     False, drop_last=True
1028 )
1029 #%% md
1030 ## Part IV: Using transformer on the toy dataset
1031 #%% md
1032 In this part, we will put all the parts together
1033 to train a transformer model. We have implemented
1034 most of the functions here for you and your task
1035 would be to use these functions to train a
1036 Transformer model. The overall tasks are divided
1037 into three parts:
1038
1039 - Implement the Transformer model using previously
1040     implemented functions
1041 - Overfitting the model
1042 - Training using complete data
1043 - Visualizing the attention
1044 #%%
1045 from tqdm import tqdm
1046 #%% md
1047 ### Implement the Transformer model
1048 #%% md
1049 The complete Transformer model combines the
1050 Encoder and Decoder stacks with input/output
1051 embeddings and a final linear layer that maps to
1052 the vocabulary size.
1053 #%% md
1054 Implement the Transformer model in `transformer.
1055 py` and run the cells below to check the final
1056 shapes of the output
1057 #%%
1058 import torch.optim as optim
1059 from transformers import Transformer
```

```

1048 inp_seq_len = 9
1049 out_seq_len = 5
1050 num_heads = 4
1051 emb_dim = 32
1052 dim_feedforward = 64
1053 dropout = 0.2
1054 num_enc_layers = 4
1055 num_dec_layers = 4
1056 vocab_len = len(vocab)
1057
1058 model = Transformer(
1059     num_heads,
1060     emb_dim,
1061     dim_feedforward,
1062     dropout,
1063     num_enc_layers,
1064     num_dec_layers,
1065     vocab_len,
1066 )
1067 for it in train_loader:
1068     it
1069     break
1070 inp, inp_pos, out, out_pos = it
1071 device = DEVICE
1072 model = model.to(device)
1073 inp_pos = inp_pos.to(device)
1074 out_pos = out_pos.to(device)
1075 out = out.to(device)
1076 inp = inp.to(device)
1077
1078
1079 model_out = model(inp.long(), inp_pos, out.long()
1080 assert model_out.size(0) == BATCH_SIZE * (
1081     out_seq_len - 1)
1082 assert model_out.size(1) == vocab_len
1083 #### Overfitting the model using small data
1084 #### md
1085 Now that we have implemented the Transformer
model, lets overfit on a small dataset. This will

```

- 1085 ensure that the implementation is correct. We keep the training and validation data same here. Before doing that, a couple of things to keep in mind:
- 1086
- 1087 - We implemented two versions of positional encodings: simple and sinusoid. For overfitting, we will use the simple positional encoding but feel free to experiment with both when training for the complete model
- 1088 - In transformers.py, we have implemented two loss functions for you. The first is the familiar cross entropy loss and second is the `LabelSmoothingLoss`. For overfitting, we will use the cross entropy loss but feel free to experiment with both while doing experiment with the complete data.
- 1089 - Usually, the training regime of Transformers start with a warmup, in other words, we train the model with a lower learning rate for some iterations and then increasing the learning rate to make the network learn faster. Intuitively, this helps you to attain a stable manifold in the loss function and then we increase the learning rate to learn faster in this stable manifold. In a way we are warming up the network to be in a stable manifold and we start training with a higher learning rate after this warm-up. For overfitting we have NOT used this warm-up as for such small data, it is okay to start with a higher learning rate but you should keep this in mind while training with the complete data. We have used two functions from a5_helper.py, `train` and `val`. Here, `train` has three parameters that you should pay attention to:
- 1090 - `warmup_interval`: Specifies the number of iterations that the network should train with a low learning rate. In other words, its the number of iterations after which the network will have the higher learning rate
- 1091 - `warmup_lr`: This is the learning rate that

```
1091 will be used during warmup.  
1092 - `lr`: This is the learning rate that will get  
    used after the warm-up. If warmup_interval is  
    None, we will start training with this learning  
    rate.  
1093  
1094 In the following cells for overfitting, we have  
    used the number of epochs as 200 but you could  
    increase this. You should get an accuracy ~1 in  
    200 epochs. It might be a little lower as well,  
    don't worry about it. It should take about a  
    minute to run the overfitting.  
1095  
1096 NOTE: When we say epoch, it means the number of  
    times we have taken a complete pass over the data  
    . One epoch typically consists of many iterations  
    that depend on the batch size.  
1097 #%%  
1098 from transformers import LabelSmoothingLoss,  
    CrossEntropyLoss  
1099 import torch.optim as optim  
1100 from transformers import Transformer  
1101 from a2_helper import train as train_transformer  
1102 from a2_helper import val as val_transformer  
1103 #%%  
1104 inp_seq_len = 9  
1105 out_seq_len = 5  
1106 num_heads = 4  
1107 emb_dim = 32  
1108 dim_feedforward = 32  
1109 dropout = 0.2  
1110 num_enc_layers = 1  
1111 num_dec_layers = 1  
1112 vocab_len = len(vocab)  
1113 BATCH_SIZE = 4  
1114 num_epochs=200 #number of epochs  
1115 lr=1e-3 #learning rate after warmup  
1116 loss_func = CrossEntropyLoss  
1117 warmup_interval = None #number of iterations for  
    warmup  
1118
```

```
1119 model = Transformer(  
1120     num_heads,  
1121     emb_dim,  
1122     dim_feedforward,  
1123     dropout,  
1124     num_enc_layers,  
1125     num_dec_layers,  
1126     vocab_len,  
1127 )  
1128 train_data = AddSubDataset(  
1129     X_train,  
1130     y_train,  
1131     convert_str_to_tokens,  
1132     SPECIAL_TOKENS,  
1133     emb_dim,  
1134     position_encoding_simple,  
1135 )  
1136 valid_data = AddSubDataset(  
1137     X_test,  
1138     y_test,  
1139     convert_str_to_tokens,  
1140     SPECIAL_TOKENS,  
1141     emb_dim,  
1142     position_encoding_simple,  
1143 )  
1144  
1145 train_loader = torch.utils.data.DataLoader(  
1146     train_data, batch_size=BATCH_SIZE, shuffle=  
     False, drop_last=True  
1147 )  
1148 valid_loader = torch.utils.data.DataLoader(  
1149     valid_data, batch_size=BATCH_SIZE, shuffle=  
     False, drop_last=True  
1150 )  
1151  
1152 small_dataset = torch.utils.data.Subset(  
1153     train_data, torch.linspace(0, len(train_data)  
         ) - 1, steps=4).long()  
1154 )  
1155 small_train_loader = torch.utils.data.DataLoader(  
1156     small_dataset, batch_size=4, pin_memory=True
```

```
1156 , num_workers=1, shuffle=False
1157 )
1158 #%%
1159
1160 #Overfitting the model
1161 trained_model = train_transformer(
1162     model,
1163     small_train_loader,
1164     small_train_loader,
1165     loss_func,
1166     num_epochs=num_epochs,
1167     lr=lr,
1168     batch_size=BATCH_SIZE,
1169     warmup_interval=warmup_interval,
1170     device=DEVICE,
1171 )
1172 #%% md
1173
1174 #%%
1175 #Overfitted accuracy
1176 print(
1177     "Overfitted accuracy: ",
1178     "{:.4f}".format(
1179         val_transformer(
1180             trained_model,
1181             small_train_loader,
1182             CrossEntropyLoss,
1183             batch_size=4,
1184             device=DEVICE,
1185             )[1]
1186     ),
1187 )
1188 #%% md
1189 ### Fitting the model using complete data
1190 #%% md
1191 Run the below cells to fit the model using the
complete data. Keep in mind the various things
you could experiment with here, losses,
positional encodings, warm up routines and
learning rates. You could also play with the size
of the model but that will require more time to
```

```
1191 train on Colab.  
1192  
1193 You should aim for final validation accuracy of ~  
     80 percent.  
1194 #%%  
1195 import torch.optim as optim  
1196 from transformers import Transformer  
1197  
1198 inp_seq_len = 9  
1199 out_seq_len = 5  
1200 BATCH_SIZE = 256  
1201  
1202 #You should change these!  
1203  
1204 num_heads = 4  
1205 emb_dim = 64  
1206 dim_feedforward = 256  
1207 dropout = 0.1  
1208 num_enc_layers = 4  
1209 num_dec_layers = 4  
1210 vocab_len = len(vocab)  
1211 loss_func = CrossEntropyLoss  
1212 poss_enc = position_encoding_sinusoid  
1213 num_epochs = 800  
1214 warmup_interval = 900  
1215 lr = 5e-4  
1216  
1217  
1218 model = Transformer(  
1219     num_heads,  
1220     emb_dim,  
1221     dim_feedforward,  
1222     dropout,  
1223     num_enc_layers,  
1224     num_dec_layers,  
1225     vocab_len,  
1226 )  
1227  
1228  
1229 train_data = AddSubDataset(  
1230     X_train,
```

```
1231     y_train,
1232     convert_str_to_tokens,
1233     SPECIAL_TOKENS,
1234     emb_dim,
1235     position_encoding_sinusoid,
1236 )
1237 valid_data = AddSubDataset(
1238     X_test,
1239     y_test,
1240     convert_str_to_tokens,
1241     SPECIAL_TOKENS,
1242     emb_dim,
1243     position_encoding_sinusoid,
1244 )
1245
1246 train_loader = torch.utils.data.DataLoader(
1247     train_data, batch_size=BATCH_SIZE, shuffle=
1248     False, drop_last=True
1249 )
1250 valid_loader = torch.utils.data.DataLoader(
1251     valid_data, batch_size=BATCH_SIZE, shuffle=
1252     False, drop_last=True
1253 )
1254 #%%
1255 #Training the model with complete data
1256 trained_model = train_transformer(
1257     model,
1258     train_loader,
1259     valid_loader,
1260     loss_func,
1261     num_epochs,
1262     lr = lr,
1263     batch_size=BATCH_SIZE,
1264     warmup_interval=warmup_interval,
1265     device=DEVICE
1266 )
1267 ARTIFACTS_DIR = ASSIGNMENT_DIR / "artifacts"
1268 ARTIFACTS_DIR.mkdir(parents=True, exist_ok=True)
1269 print("Using ARTIFACTS_DIR:", ARTIFACTS_DIR)
1270 weights_path = str(ARTIFACTS_DIR / "transformer.
1271 pt")
```

```
1269 torch.save(trained_model.state_dict(),
   weights_path)
1270 #%% md
1271 Run the cell below to get the accuracy on the
 validation dataset.
1272 #%%
1273 #Final validation accuracy
1274 print(
1275     "Final Model accuracy: ",
1276     "{:.4f}".format(
1277         val_transformer(
1278             trained_model, valid_loader,
1279             LabelSmoothingLoss, 4, device=DEVICE
1280         )[1]
1281     ),
1282 #%% md
1283 ## Visualize and Inference: Model in action
1284 #%% md
1285 Now that we have trained a model, lets look at
 the final results. We will first look at the
 results from the validation data and visualize
 the attention weights (remember the self.
 weights_softmax?). These attention weights should
 give you some intuition about what the network
 learns. We have implemented everything for you
 here and the intention is to help you probe the
 model and understand about what does the network
 learn.
1286 #%%
1287 import seaborn
1288 from a2_helper import inference
1289 #%% md
1290 ### Results from the validation data
1291
1292 In the below cell we pick the very first data
 point in the validation data and find the result
 on it
1293 #%%
1294 for it in valid_loader:
1295     it
```

```
1296     break
1297 inp, inp_pos, out, out_pos = it
1298 opposite_tokens_to_str = {v: k for k, v in
1299     convert_str_to_tokens.items()}
1300 device = torch.device("mps")
1300 model = model.to(device)
1301 inp_pos = inp_pos.to(device)
1302 out_pos = out_pos.to(device)
1303 out = out.to(device)
1304 inp = inp.to(device)
1305
1306 inp_exp = inp[:1, :]
1307 inp_exp_pos = inp_pos[:1]
1308 out_pos_exp = out_pos[:1, :]
1309 inp_seq = [opposite_tokens_to_str[w.item()] for w
1310     in inp_exp[0]]
1310 print(
1311     "Input sequence: \n",
1312     inp_seq[0]
1313     + " "
1314     + inp_seq[1]
1315     + " "
1316     + inp_seq[2]
1317     + inp_seq[3]
1318     + " "
1319     + inp_seq[4]
1320     + " "
1321     + inp_seq[5]
1322     + " "
1323     + inp_seq[6]
1324     + inp_seq[7]
1325     + " "
1326     + inp_seq[8],
1327 )
1328 #%%
1329 out_seq_ans, _ = inference(
1330     trained_model, inp_exp, inp_exp_pos,
1331     out_pos_exp, out_seq_len
1331 )
1332
1333 trained_model.eval()
```

```
1334
1335 print("Output Sequence:", end="\t")
1336 res = "BOS "
1337 for i in range(1, out_seq_ans.size(1)):
1338     sym = opposite_tokens_to_str[out_seq_ans[0, i
1339 ].item()]
1340     if sym == "EOS":
1341         break
1342     res += sym + " "
1343 print(res)
1344 #%% md
1345 ### Pick your own probing example
1346 In the cell below, you could feed in an example
in the input style, changing the variable
`custom_seq`. We have filled a placeholder
expression for you, but feel free to change it.
1347 #%%
1348 custom_seq = "BOS POSITIVE 02 subtract POSITIVE
15 EOS"
1349 #%%
1350 out = preprocess_input_sequence(custom_seq,
convert_str_to_tokens, SPECIAL_TOKENS)
1351 inp_exp = torch.tensor(out).to(DEVICE)
1352
1353 out_seq_ans, model_for_visv = inference(
1354     trained_model, inp_exp, inp_exp_pos,
out_pos_exp, out_seq_len
1355 )
1356
1357 trained_model.eval()
1358
1359 print("Output Sequence:", end="\t")
1360 res = "BOS "
1361 for i in range(1, out_seq_ans.size(1)):
1362     sym = opposite_tokens_to_str[out_seq_ans[0, i
1363 ].item()]
1364     if sym == "EOS":
1365         break
1366     res += sym + " "
1367 print(res)
```

```
1367 #%% md
1368 ### Visualize the attention weights
1369
1370 In this part we will visualize the attention
   weights for the specific custom input you fed as
   input. There are separate heatmaps for encoder
   and the decoder. The lighter value in color shows
   higher association between the token present in
   that row and column, and darker color shows a
   weak relation between them
1371 #%%
1372 from a2_helper import draw
1373 import seaborn
1374 #%%
1375 target_exp = res.split()
1376 #%%
1377 for layer in range(num_enc_layers):
1378     fig, axs = plt.subplots(1, num_heads, figsize
1379     =(20, 10))
1380     print("Encoder Block Number", layer + 1)
1381     for h in range(num_heads):
1382         draw(
1383             trained_model.encoder.layers[layer]
1384             .MultiHeadBlock.heads[h]
1385             .weights_softmax.data.cpu()
1386             .numpy()[0],
1387             inp_seq,
1388             inp_seq if h == 0 else [],
1389             ax=axs[h],
1390             )
1391     plt.show()
1392 plt.close()
1393 #%%
1394 for layer in range(num_dec_layers):
1395     fig, axs = plt.subplots(1, num_heads, figsize
1396     =(20, 10))
1397     print("Decoder Block number ", layer + 1)
1398     print("Decoder Self Attention", layer + 1)
1399     for h in range(num_heads):
```

```
1400     draw(
1401         trained_model.decoder.layers[layer]
1402         .attention_self.heads[h]
1403         .weights_softmax.data.cpu()
1404         .numpy()[0],
1405         target_exp,
1406         target_exp if h == 0 else [],
1407         ax=axs[h],
1408     )
1409     plt.show()
1410     print("Decoder Cross attention", layer + 1)
1411     fig, axs = plt.subplots(1, 4, figsize=(20, 10))
1412     for h in range(num_heads):
1413         draw(
1414             trained_model.decoder.layers[layer]
1415             .attention_cross.heads[h]
1416             .weights_softmax.data.cpu()
1417             .numpy()[0],
1418             inp_seq,
1419             target_exp if h == 0 else [],
1420             ax=axs[h],
1421         )
1422     plt.show()
1423 %% md
1424 # Submit Your Work
1425 After completing both notebooks for this assignment (`Transformers.ipynb` and `rnn_lstm_captioning.ipynb`), run the following cell to create a `*.zip` file for you to download and turn in.
1426
1427 **Please MANUALLY SAVE every `*.ipynb` and `*.py` files before executing the following cell:**
```

```
1428 %% md
1429 from dl.submit import make_assignment2_submission
1430
1431 make_assignment2_submission("./")
```

```
1 #%%
2 # This mounts your Google Drive to the Colab VM.
3 from google.colab import drive
4 drive.mount('/content/drive')
5
6 # TODO: Enter the foldername in your Drive where
7 # you have saved the unzipped
8 # assignment folder, e.g. 'dl/assignments/
9 # assignment2'
10 FOLDERNAME = "dl/assignments/assignment2/"
11 assert FOLDERNAME is not None, "[!] Enter the
12 # foldername."
13
14 import sys
15 sys.path.append('/content/drive/My Drive/{}'.format
16 (FOLDERNAME))
17
18 # Change to the assignment directory
19 %cd /content/drive/My\ Drive/$FOLDERNAME
20 ## Using GPU
21
22 Go to `Runtime > Change runtime type` and set `Hardware accelerator` to `GPU`. This will reset Colab. **Rerun the top cell to mount your Drive again.**
23
24 # Self-Supervised Learning
25
26 ## What is self-supervised learning?
27 Modern day machine learning requires lots of labeled data. But often times it's challenging and/or expensive to obtain large amounts of human-labeled data. Is there a way we could ask machines to automatically learn a model which can generate good visual representations without a labeled dataset? Yes, enter self-supervised learning!
```

28

29 Self-supervised learning (SSL) allows models to automatically learn a "good" representation space using the data in a given dataset without the need for their labels. Specifically, if our dataset were a bunch of images, then self-supervised learning allows a model to learn and generate a "good" representation vector for images.

30

31 The reason SSL methods have seen a surge in popularity is because the learnt model continues to perform well on other datasets as well i.e. new datasets on which the model was not trained on!

32

33 ## What makes a "good" representation?

34 A "good" representation vector needs to capture the important features of the image as it relates to the rest of the dataset. This means that images in the dataset representing semantically similar entities should have similar representation vectors , and different images in the dataset should have different representation vectors. For example, two images of an apple should have similar representation vectors, while an image of an apple and an image of a banana should have different representation vectors.

35

36 ## Contrastive Learning: SimCLR

37 Recently, [SimCLR](<https://arxiv.org/pdf/2002.05709.pdf>) introduces a new architecture which uses ** contrastive learning** to learn good visual representations. Contrastive learning aims to learn similar representations for similar images and different representations for different images. As we will see in this notebook, this simple idea allows us to train a surprisingly good model without using any labels.

38

39 Specifically, for each image in the dataset, SimCLR generates two differently augmented views of that image, called a **positive pair**. Then, the model

```

39 is encouraged to generate similar representation
vectors for this pair of images. See below for an
illustration of the architecture (Figure 2 from the
paper).


---


40 #%%
41 # Run this cell to view the SimCLR architecture.
42 from IPython.display import Image
43 Image('images/simclr_fig2.png', width=500)


---


44 #%% md
45 Given an image  $\mathbf{x}$ , SimCLR uses two different
data augmentation schemes  $\mathbf{t}$  and  $\mathbf{t}'$  to
generate the positive pair of images  $\tilde{\mathbf{x}}$ 
 $\mathbf{i}$  and  $\tilde{\mathbf{x}}_j$ .  $f$  is a basic encoder
net that extracts representation vectors from the
augmented data samples, which yields  $\mathbf{h}_i$  and
 $\mathbf{h}_j$ , respectively. Finally, a small neural
network projection head  $g$  maps the representation
vectors to the space where the contrastive loss is
applied. The goal of the contrastive loss is to
maximize agreement between the final vectors  $\mathbf{z}_i$ 
 $= g(\mathbf{h}_i)$  and  $\mathbf{z}_j = g(\mathbf{h}_j)$ . We will
discuss the contrastive loss in more detail later,
and you will get to implement it.


---


46
47 After training is completed, we throw away the
projection head  $g$  and only use  $f$  and the
representation  $\mathbf{h}$  to perform downstream tasks,
such as classification. You will get a chance to
finetune a layer on top of a trained SimCLR model
for a classification task and compare its
performance with a baseline model (without self-
supervised learning).


---


48 #%% md
49 ## Pretrained Weights
50 For your convenience, we have given you pretrained
weights (trained for ~18 hours on CIFAR-10) for the
SimCLR model. Run the following cell to download
pretrained model weights to be used later. (This
will take ~1 minute)


---


51 #%%
52 %%bash

```

```
53 DIR=pretrained_model/
54 if [ ! -d "$DIR" ]; then
55     mkdir "$DIR"
56 fi
57
58 URL=http://downloads.cs.stanford.edu/downloads/
      cs231n/pretrained_simclr_model.pth
59 # Try this if above doesn't work.
60 # URL=http://cs231n.stanford.edu/2025/storage/a3/
      pretrained_simclr_model.pth
61 FILE=pretrained_model/pretrained_simclr_model.pth
62 if [ ! -f "$FILE" ]; then
63     echo "Downloading weights..."
64     wget "$URL" -O "$FILE"
65 fi
66 #%%
67 # Setup cell.
68 %pip install thop
69 import torch
70 import os
71 import importlib
72 import pandas as pd
73 import numpy as np
74 import torch.optim as optim
75 import torch.nn as nn
76 import random
77 from thop import profile, clever_format
78 from torch.utils.data import DataLoader
79 from torchvision.datasets import CIFAR10
80 import matplotlib.pyplot as plt
81 %matplotlib inline
82
83 # load_ext autoreload
84 # autoreload 2
85
86 #device = torch.device("cuda" if torch.cuda.
     is_available() else "cpu")
87 # for mac with m3 chip use:
88 device = torch.device("mps" if torch.backends.mps.
     is_available() else "cpu")
89 print(f"Using device: {device}")
```

```

90 #%% md
91 # Data Augmentation
92
93 Our first step is to perform data augmentation.
  Implement the `compute_train_transform()` function
    in `dl/simclr/data_utils.py` to apply the
      following random transformations:
94
95 1. Randomly resize and crop to 32x32.
96 2. Horizontally flip the image with probability 0.
  5
97 3. With a probability of 0.8, apply color jitter (
    see `compute_train_transform()` for definition)
98 4. With a probability of 0.2, convert the image to
    grayscale
99 #%% md
100 Now complete `compute_train_transform()` and `CIFAR10Pair.__getitem__()` in `dl/simclr/
  data_utils.py` to apply the data augmentation
  transform and generate **$\tilde{x}_i$** and **$\tilde{x}_j$**.
101
102
103 #%% md
104 Test to make sure that your data augmentation code
  is correct:
105 #%%
106 from dl.simclr.data_utils import *
107 from dl.simclr.contrastive_loss import *
108
109 answers = torch.load('simclr_sanity_check.key')
110 #%% md
111 # Base Encoder and Projection Head
112 The next steps are to apply the base encoder and
  projection head to the augmented samples **$\tilde{x}_i$** and **$\tilde{x}_j$**.
113
114 The base encoder $f$ extracts representation
  vectors for the augmented samples. The SimCLR
  paper found that using deeper and wider models
  improved performance and thus chose [ResNet](https

```

```

114 ://arxiv.org/pdf/1512.03385.pdf) to use as the
base encoder. The output of the base encoder are
the representation vectors **$h_i = f(\tilde{x}_i$**
)** and **$h_j = f(\tilde{x}_j$)**.
115
116 The projection head $g$ is a small neural network
that maps the representation vectors **$h_i$** and
**$h_j$** to the space where the contrastive loss
is applied. The paper found that using a
nonlinear projection head improved the
representation quality of the layer before it.
Specifically, they used a MLP with one hidden
layer as the projection head $g$. The contrastive
loss is then computed based on the outputs **$z_i$**
= $g(h_i)$ and **$z_j$** = $g(h_j)$.
117
118 We provide implementations of these two parts in `dl/simclr/model.py`. Please skim through the file
and make sure you understand the implementation.


---


119 #%% md
120 # SimCLR: Contrastive Loss
121
122 A mini-batch of $N$ training images yields a total
of $2N$ data-augmented examples. For each
positive pair $(i, j)$ of augmented examples, the
contrastive loss function aims to maximize the
agreement of vectors $z_i$ and $z_j$. Specifically,
the loss is the normalized temperature-scaled
cross entropy loss and aims to maximize the
agreement of $z_i$ and $z_j$ relative to all other
augmented examples in the batch:


---


123 #%% md
124 $$
125 l \; (i, j) = -\log \frac{\exp (\text{sim}(z_i, z_j) / \tau)}{\sum_{k=1}^{2N} \mathbb{1}_{k \neq i} \exp (\text{sim}(z_i, z_k) / \tau)}
126 $$


---


127 #%% md
128 where $\mathbb{1} \in \{0, 1\}$ is an indicator
function that outputs $1$ if $k \neq i$ and $0$ otherwise. $\tau$ is a temperature parameter that

```

```

128 determines how fast the exponentials increase.
129
130 sim$(z_i, z_j) = \frac{z_i \cdot z_j}{\|z_i\| \|z_j\|} is the (normalized) dot product
between vectors $z_i$ and $z_j$. The higher the
similarity between $z_i$ and $z_j$, the larger the
dot product is, and the larger the numerator
becomes. The denominator normalizes the value by
summing across $z_i$ and all other augmented
examples $k$ in the batch. The range of the
normalized value is $(0, 1)$, where a high score
close to $1$ corresponds to a high similarity
between the positive pair $(i, j)$ and low
similarity between $i$ and other augmented
examples $k$ in the batch. The negative log then
maps the range $(0, 1)$ to the loss values $(\inf,
0)$.
131
132 The total loss is computed across all positive
pairs $(i, j)$ in the batch. Let $z = [z_1, z_2
, \dots, z_{2N}]$ include all the augmented examples
in the batch, where $z_1 \dots z_N$ are outputs
of the left branch, and $z_{N+1} \dots z_{2N}$ are
outputs of the right branch. Thus, the positive
pairs are $(z_k, z_{k+N})$ for $\forall k \in
[1, N]$.
133
134 Then, the total loss $L$ is:
135 %% md
136 $$
137 L = \frac{1}{2N} \sum_{k=1}^N [ l(k, ; k+N) + l
(k+N, ; k) ; ]
138 $$
139 %% md
140 **NOTE:** this equation is slightly different from
the one in the paper. We've rearranged the
ordering of the positive pairs in the batch, so
the indices are different. The rearrangement makes
it easier to implement the code in vectorized
form.
141

```

```

142 We'll walk through the steps of implementing the
    loss function in vectorized form. Implement the
    functions `sim`, `simclr_loss_naive` in `dl/simclr
    /contrastive_loss.py`. Test your code by running
    the sanity checks below.
143 #%%
144 from dl.simclr.contrastive_loss import *
145 answers = torch.load('simclr_sanity_check.key')
146 #%%
147 def test_sim(left_vec, right_vec, correct_output):
148     output = sim(left_vec, right_vec).cpu().numpy()
149     print("Maximum error in sim: %g"%rel_error(
150         correct_output.numpy(), output))
151 # Should be less than 1e-07.
152 test_sim(answers['left'][0], answers['right'][0],
153 answers['sim'][0])
154 test_sim(answers['left'][1], answers['right'][1],
155 answers['sim'][1])
156 #%%
157 def test_loss_naive(left, right, tau,
158 correct_output):
159     naive_loss = simclr_loss_naive(left, right,
160     tau).item()
161     print("Maximum error in simclr_loss_naive: %g"
162         "%rel_error(correct_output, naive_loss))")
163 # Should be less than 1e-07.
164 test_loss_naive(answers['left'], answers['right'],
165 [5.0, answers['loss']['5.0']])
166 test_loss_naive(answers['left'], answers['right'],
167 [1.0, answers['loss']['1.0']])
168 #%% md
169 Now implement the vectorized version by
    implementing `sim_positive_pairs`,
    `compute_sim_matrix`, `simclr_loss_vectorized` in
    `dl/simclr/contrastive_loss.py`. Test your code by
    running the sanity checks below.
170 #%%
171 def test_sim_positive_pairs(left, right,

```

```

165 correct_output):
166     sim_pair = sim_positive_pairs(left, right).cpu()
167     print("Maximum error in sim_positive_pairs: %g"
168           "%rel_error(correct_output.numpy(), sim_pair)")
169 # Should be less than 1e-07.
170 test_sim_positive_pairs(answers['left'], answers['right'],
171                         answers['sim'])
171 #%%
172 def test_sim_matrix(left, right, correct_output):
173     out = torch.cat([left, right], dim=0)
174     sim_matrix = compute_sim_matrix(out).cpu()
175     assert torch.isclose(sim_matrix,
176                          correct_output).all(), "correct: {}. got: {}".
177                          format(correct_output, sim_matrix)
176     print("Test passed!")
177
178 test_sim_matrix(answers['left'], answers['right'],
179                  answers['sim_matrix'])
179 #%%
180 def test_loss_vectorized(left, right, tau,
181                         correct_output):
181     vec_loss = simclr_loss_vectorized(left, right,
182                                         tau, device=left.device).item()
182     print("Maximum error in loss_vectorized: %g"%
183           "%rel_error(correct_output, vec_loss)")
183
184 # Should be less than 1e-07.
185 test_loss_vectorized(answers['left'], answers['right'],
186                      5.0, answers['loss']['5.0'])
186 test_loss_vectorized(answers['left'], answers['right'],
187                      1.0, answers['loss']['1.0'])
187 #%% md
188 # Implement the train function
189 Complete the `train()` function in `dl/simclr/
190 utils.py` to obtain the model's output and use
191 `simclr_loss_vectorized` to compute the loss. (
192 Please take a look at the `Model` class in `dl/
193 simclr/model.py` to understand the model pipeline
194 and the returned values)

```

```
190 #%%
191 from dl.simclr.data_utils import *
192 from dl.simclr.model import *
193 from dl.simclr.utils import *
194 #%% md
195 ### Train the SimCLR model
196
197 Run the following cells to load in the pretrained
weights and continue to train a little bit more.
This part will take ~10 minutes and will output to
`pretrained_model/trained_simclr_model.pth`.
198
199 **NOTE:** Don't worry about logs such as '_[WARN]
Cannot find rule for ...'. These are related to
another module used in the notebook. You can
verify the integrity of your code changes through
our provided prompts and comments.
200 #%%
201 # Do not modify this cell.
202 feature_dim = 128
203 temperature = 0.5
204 k = 200
205 batch_size = 64
206 epochs = 1
207 temperature = 0.5
208 percentage = 0.5
209 pretrained_path = './pretrained_model/
    pretrained_simclr_model.pth'
210
211 # Prepare the data.
212 train_transform = compute_train_transform()
213 train_data = CIFAR10Pair(root='data', train=True,
    transform=train_transform, download=True)
214 train_data = torch.utils.data.Subset(train_data,
    list(np.arange(int(len(train_data)*percentage))))
215 train_loader = DataLoader(train_data, batch_size=
    batch_size, shuffle=True, num_workers=16,
    pin_memory=True, drop_last=True)
216 test_transform = compute_test_transform()
217 memory_data = CIFAR10Pair(root='data', train=True
    , transform=test_transform, download=True)
```

```

218 memory_loader = DataLoader(memory_data, batch_size
                                =batch_size, shuffle=False, num_workers=16,
                                pin_memory=True)
219 test_data = CIFAR10Pair(root='data', train=False,
                           transform=test_transform, download=True)
220 test_loader = DataLoader(test_data, batch_size=
                                batch_size, shuffle=False, num_workers=16,
                                pin_memory=True)
221
222 # Set up the model and optimizer config.
223 model = Model(feature_dim)
224 model.load_state_dict(torch.load(pretrained_path,
                                    map_location='cpu'), strict=False)
225 model = model.to(device)
226 flops, params = profile(model, inputs=(torch.randn
(1, 3, 32, 32).to(device),))
227 flops, params = clever_format([flops, params])
228 print('# Model Params: {} FLOPs: {}'.format(params
, flops))
229 optimizer = optim.Adam(model.parameters(), lr=1e-3
, weight_decay=1e-6)
230 c = len(memory_data.classes)
231
232 # Training loop.
233 results = {'train_loss': [], 'test_acc@1': [], 'test_acc@5': []} #<< -- output
234
235 if not os.path.exists('results'):
236     os.mkdir('results')
237 best_acc = 0.0
238 for epoch in range(1, epochs + 1):
239     train_loss = train(model, train_loader,
                           optimizer, epoch, epochs, batch_size=batch_size,
                           temperature=temperature, device=device)
240     results['train_loss'].append(train_loss)
241     test_acc_1, test_acc_5 = test(model,
                           memory_loader, test_loader, epoch, epochs, c, k=k
                           , temperature=temperature, device=device)
242     results['test_acc@1'].append(test_acc_1)
243     results['test_acc@5'].append(test_acc_5)
244

```

```

245      # Save statistics.
246      if test_acc_1 > best_acc:
247          best_acc = test_acc_1
248          torch.save(model.state_dict(), './
249          pretrained_model/trained_simclr_model.pth)'
250 ##% md
251 # Finetune a Linear Layer for Classification!
252 Now it's time to put the representation vectors to
the test!
253
254 We remove the projection head from the SimCLR
model and slap on a linear layer to finetune for a
simple classification task. All layers before the
linear layer are frozen, and only the weights in
the final linear layer are trained. We compare the
performance of the SimCLR + finetuning model
against a baseline model, where no self-supervised
learning is done beforehand, and all weights in
the model are trained. You will get to see for
yourself the power of self-supervised learning and
how the learned representation vectors improve
downstream task performance.
255
256 ##% md
257 ## Baseline: Without Self-Supervised Learning
258 First, let's take a look at the baseline model. We
'll remove the projection head from the SimCLR
model and slap on a linear layer to finetune for a
simple classification task. No self-supervised
learning is done beforehand, and all weights in
the model are trained. Run the following cells.
259
260 **NOTE:** Don't worry if you see low but
reasonable performance.
261 ##%
262 class Classifier(nn.Module):
263     def __init__(self, num_class):
264         super(Classifier, self).__init__()
265
266         # Encoder.

```

```
267         self.f = Model().f
268
269         # Classifier.
270         self.fc = nn.Linear(2048, num_class, bias=
True)
271
272     def forward(self, x):
273         x = self.f(x)
274         feature = torch.flatten(x, start_dim=1)
275         out = self.fc(feature)
276         return out
277 #%%
278 # Do not modify this cell.
279 feature_dim = 128
280 temperature = 0.5
281 k = 200
282 batch_size = 128
283 epochs = 10
284 percentage = 0.1
285
286 train_transform = compute_train_transform()
287 train_data = CIFAR10(root='data', train=True,
transform=train_transform, download=True)
288 trainset = torch.utils.data.Subset(train_data,
list(np.arange(int(len(train_data)*percentage))))
289 train_loader = DataLoader(trainset, batch_size=
batch_size, shuffle=True, num_workers=16,
pin_memory=True)
290 test_transform = compute_test_transform()
291 test_data = CIFAR10(root='data', train=False,
transform=test_transform, download=True)
292 test_loader = DataLoader(test_data, batch_size=
batch_size, shuffle=False, num_workers=16,
pin_memory=True)
293
294 model = Classifier(num_class=len(train_data.
classes)).to(device)
295 for param in model.f.parameters():
296     param.requires_grad = False
297
298 flops, params = profile(model, inputs=(torch.randn
```

```

298 (1, 3, 32, 32).to(device),))
299 flops, params = clever_format([flops, params])
300 print('# Model Params: {} FLOPs: {}'.format(params,
301 , flops))
301 optimizer = optim.Adam(model.fc.parameters(), lr=
302 1e-3, weight_decay=1e-6)
302 no_pretrain_results = {'train_loss': [], 'train_acc@1': [],
303 'test_loss': [], 'test_acc@1': [],
304 'test_acc@5': []}
304
305 best_acc = 0.0
306 for epoch in range(1, epochs + 1):
307     train_loss, train_acc_1, train_acc_5 =
308         train_val(model, train_loader, optimizer, epoch,
309         epochs, device=device)
310     no_pretrain_results['train_loss'].append(
311         train_loss)
312     no_pretrain_results['train_acc@1'].append(
313         train_acc_1)
314     no_pretrain_results['train_acc@5'].append(
315         train_acc_5)
316     test_loss, test_acc_1, test_acc_5 = train_val(
317         model, test_loader, None, epoch, epochs, device=
318         device)
318     no_pretrain_results['test_loss'].append(
319         test_loss)
320     no_pretrain_results['test_acc@1'].append(
321         test_acc_1)
322     no_pretrain_results['test_acc@5'].append(
323         test_acc_5)
323
324     if test_acc_1 > best_acc:
325         best_acc = test_acc_1
326
327
328 # Print the best test accuracy.
329 print('Best top-1 accuracy without self-supervised
330       learning: ', best_acc)
331 ## With Self-Supervised Learning
332
333 Let's see how much improvement we get with self-

```

```
323 supervised learning. Here, we pretrain the SimCLR
model using the simclr loss you wrote, remove the
projection head from the SimCLR model, and use a
linear layer to finetune for a simple
classification task.
324 #%%
325 # Do not modify this cell.
326 feature_dim = 128
327 temperature = 0.5
328 k = 200
329 batch_size = 128
330 epochs = 10
331 percentage = 0.1
332 pretrained_path = './pretrained_model/
trained_simclr_model.pth'
333
334 train_transform = compute_train_transform()
335 train_data = CIFAR10(root='data', train=True,
transform=train_transform, download=True)
336 trainset = torch.utils.data.Subset(train_data,
list(np.arange(int(len(train_data)*percentage))))
337 train_loader = DataLoader(trainset, batch_size=
batch_size, shuffle=True, num_workers=16,
pin_memory=True)
338 test_transform = compute_test_transform()
339 test_data = CIFAR10(root='data', train=False,
transform=test_transform, download=True)
340 test_loader = DataLoader(test_data, batch_size=
batch_size, shuffle=False, num_workers=16,
pin_memory=True)
341
342 model = Classifier(num_class=len(train_data.
classes))
343 model.load_state_dict(torch.load(pretrained_path,
map_location='cpu'), strict=False)
344 model = model.to(device)
345 for param in model.parameters():
346     param.requires_grad = False
347
348 flops, params = profile(model, inputs=(torch.randn
(1, 3, 32, 32).to(device),))
```

```

349 flops, params = clever_format([flops, params])
350 print('# Model Params: {} FLOPs: {}'.format(params,
351 , flops))
351 optimizer = optim.Adam(model.fc.parameters(), lr=
352 1e-3, weight_decay=1e-6)
352 pretrain_results = {'train_loss': [], 'train_acc@1':
353 ': [], 'train_acc@5': [],
353 'test_loss': [], 'test_acc@1': [], 'test_acc@5': []}
354
355 best_acc = 0.0
356 for epoch in range(1, epochs + 1):
357     train_loss, train_acc_1, train_acc_5 =
357      train_val(model, train_loader, optimizer, epoch,
357      epochs, device=device)
358     pretrain_results['train_loss'].append(
358      train_loss)
359     pretrain_results['train_acc@1'].append(
359      train_acc_1)
360     pretrain_results['train_acc@5'].append(
360      train_acc_5)
361     test_loss, test_acc_1, test_acc_5 = train_val(
361      model, test_loader, None, epoch, epochs, device=
361      device)
362     pretrain_results['test_loss'].append(test_loss
362 )
363     pretrain_results['test_acc@1'].append(
363      test_acc_1)
364     pretrain_results['test_acc@5'].append(
364      test_acc_5)
365     if test_acc_1 > best_acc:
366         best_acc = test_acc_1
367
368 # Print the best test accuracy. You should see a
368 best top-1 accuracy of >=70%.
369 print('Best top-1 accuracy with self-supervised
369 learning: ', best_acc)
370 #%% md
371 ### Plot your Comparison
372
373 Plot the test accuracies between the baseline

```

```
373 model (no pretraining) and same model pretrained  
with self-supervised learning.

---

374 #%%  
375 plt.plot(no_pretrain_results['test_acc@1'], label=  
"Without Pretrain")  
376 plt.plot(pretrain_results['test_acc@1'], label="  
With Pretrain")  
377 plt.xlabel('Epochs')  
378 plt.ylabel('Accuracy')  
379 plt.title('Test Top-1 Accuracy')  
380 plt.legend()  
381 plt.show()

---

382 #%%  
383 # Note: Run the submission cell in CLIP_DINO.ipynb  
(the final notebook) to generate submission files  
.
```

```
1 #%%
2 # This mounts your Google Drive to the Colab VM.
3 from google.colab import drive
4 drive.mount('/content/drive', force_remount=True)
5
6 # TODO: Enter the foldername in your Drive where
7 # you have saved the unzipped
8 # assignment folder, e.g. 'dl/assignments/
9 # assignment2/'
10 FOLDERNAME = "dl/assignments/assignment2/"
11 assert FOLDERNAME is not None, "[!] Enter the
12 # foldername."
13
14 import sys
15 sys.path.append('/content/drive/My Drive/{}'.format
16 (FOLDERNAME))
17
18 #%%
19 # This downloads the COCO dataset to your Drive if
20 # it doesn't already exist
21 # (you should already have this dataset from a
22 # previous notebook!)
23 # Uncomment the following if you don't have it.
24 # %cd /content/drive/My\ Drive/$FOLDERNAME/dl/
25 # datasets/
26 # !bash get_coco_captioning.sh
27 # %cd /content/drive/My\ Drive/$FOLDERNAME
28 #%%
29 # cd /Users/jeremy/Desktop/deep - Idan/exercise/ex2/
30 # DL-Attention/
31 #%%
```

```
29 # Some useful python libraries
30 ! pip install ftfy regex tqdm
31 ! pip install git+https://github.com/openai/CLIP.
32 git
```

```
32 ! pip install decord
33 #%% md
34 # State-of-the-Art Pretrained Image Models
35
36 In the previous exercise, you learned about [SimCLR]
[https://arxiv.org/abs/2002.05709] and how
contrastive self-supervised learning can be used to
learn meaningful image representations. In this
notebook, we will explore two more recent models
that also aim to learn high-quality visual
representations and have demonstrated strong and
robust performance on a variety of downstream tasks
.
37
38
39 First, we will examine the [CLIP](https://github.com/openai/CLIP) model. Like SimCLR, CLIP uses a
contrastive learning objective, but instead of
contrasting two augmented views of the same image,
it contrasts two different modalities: text and
image. To train CLIP, OpenAI collected a large
dataset of ~400M image-text pairs from the internet
, including sources like Wikipedia and image alt
text. The resulting model learns rich, high-level
image features and has achieved impressive zero-
shot performance on many vision benchmarks.
40
41 Next, we will explore [DINO](https://github.com/facebookresearch/dino), a self-supervised learning
method for vision tasks that applies contrastive
learning in a self-distillation framework with
multi-crop augmentation strategy. The authors
showed that the features learned by DINO ViTs are
fine-grained and semantically rich with explicit
information about the semantic segmentation of the
image.
42
43
44
45 #%% md
46 # CLIP
```

47

- 48 As explained above, CLIP's training objective incorporates both text and images, building upon the principles of contrastive learning. Consider this quote from the SimCLR notebook:
- 49 >The goal of the contrastive loss is to maximize agreement between the final vectors $z_i = g(h_i)$ and $z_j = g(h_j)$.

50

- 51 Similarly, CLIP is trained to maximize agreement between two vectors. However, because these vectors come from different modalities, CLIP uses two separate encoders: a transformer-based Text Encoder and a Vision Transformer (ViT)-based Image Encoder. Note that some smaller, more efficient versions of CLIP use a ResNet as the Image Encoder instead of a ViT.

52

- 53 Run the cell below to visualize the training and inference pipeline of CLIP.

54

- 55 During the pretraining phase, each batch consists of multiple images along with their corresponding captions. Each image is independently processed by an Image Encoder—typically a visual model like a Vision Transformer (ViT) or a Convolutional Neural Network (ConvNet)—which produces an image embedding I_n . Likewise, each caption is independently processed by a Text Encoder to generate a corresponding text embedding T_n . Next, we compute the pairwise similarities between all image-text combinations, meaning each image is compared with every caption, and vice versa. The training objective is to maximize the similarity scores along the diagonal of the resulting similarity matrix -- that is, the scores for the matching image-caption pairs (I_n, T_n) . Through backpropagation, the model learns to assign higher similarity scores to true matches than to mismatched pairs.

56

```
57 Through this setup, CLIP effectively learns to
    represent images and texts in a shared latent space
    . In this space, semantic concepts are encoded in a
    modality-independent way, enabling meaningful
    cross-modal comparisons between visual and textual
    inputs.
58
59
60 #%%
61 from IPython.display import Image as ColabImage
62 ColabImage(f'CLIP.png')
63 #%% md
64 **Inline Question 1** -
65
66 Why does CLIP's learning depend on the batch size?
    If the batch size is fixed, what strategy can we
    use to learn rich image features?
67
68 $\color{blue}{\textit{Your Answer:}}$%
69
70 CLIP's learning depends on batch size because of
    the nature of its contrastive loss.
71 In each training batch, every image-text pair
    treats all other samples in the batch as negative
    examples. Therefore, increasing the batch size
    increases the number and diversity of negatives
    seen at each training step. This creates stronger
    contrastive pressure, encouraging the model to
    learn more fine-grained, discriminative, and
    semantically meaningful representations, which
    generally improves generalization.
72
73 However, batch size is often limited by memory
    constraints. When the batch size must remain fixed
    , we can still obtain rich representations by
    increasing the number of effective negatives
    without increasing the batch size.
74
75 A common strategy is MoCo-style memory queues:
76
77 - Only the current batch participates in gradient
```

```
77 computation.  
78  
79 - Embeddings from previous batches are stored in a  
    queue and reused as additional negative examples.  
80  
81 - These queued embeddings are treated as constants  
    (no gradients flow through them).  
82  
83 This dramatically increases the number of  
    negatives available for contrastive learning  
    without increasing memory usage for  
    backpropagation, effectively decoupling batch size  
    from the number of negatives.  
84  
85 In addition, other complementary techniques can  
    further enrich learning:  
86  
87 - Hard or confusing negative mining, which focuses  
    on negatives that are close to the positives in  
    embedding space.  
88  
89 - Strong data augmentation, which increases view  
    diversity and encourages invariance.  
90  
91 - Cross-batch or distributed negatives, where  
    embeddings are shared across devices or steps.  
92  
93  
94  
95

---

  
96 #%% md  
97 # Loading COCO dataset  
98  
99 We'll use the same captioning dataset you used to  
    train your RNN captioning model, but instead of  
    generating the captions lets see if we can match  
    each image to the correct caption.  
100 #%%  
101 # load_ext autoreload  
102 # autoreload 2  
103
```

```
104 import time, os, json
105 import numpy as np
106 import matplotlib.pyplot as plt
107 import torch
108 import clip
109 import torch
110 from tqdm.auto import tqdm
111
112 from PIL import Image
113 from dl.clip_dino import *
114
115 def rel_error(x, y):
116     """Returns relative error."""
117     return np.max(np.abs(x - y) / (np.maximum(1e-
118                                     10, np.abs(x) + np.abs(y))))
119 #%%
120 from dl.coco_utils import load_coco_data,
121     sample_coco_minibatch, decode_captions
121 from dl.image_utils import image_from_url
122 #%%
123 # Load COCO data from disk into a dictionary.
124 # this is the same dataset you used for the RNN
125 captioning notebook :)
125 data = load_coco_data(pca_features=True)
126
127 # Print out all the keys and values from the data
128 # dictionary.
128 for k, v in data.items():
129     if type(v) == np.ndarray:
130         print(k, type(v), v.shape, v.dtype)
131     else:
132         print(k, type(v), len(v))
133 #%%
134 # we're just using the loaded captions from COCO,
135 # so we need to decode them and get rid of the
136 # special tokens.
135 decoded_captions= []
136 for caption in data['val_captions']:
137     caption = decode_captions(caption, data['
137 idx_to_word'])\
```

```

138     .replace('<START>', '') \
139     .replace('<END>', '') \
140     .replace('<UNK>', '') \
141     .strip()
142     decoded_captions.append(caption)
143 #%%
144 # lets get 10 examples
145 mask = np.array([135428, 122586, 122814, 133173,
176639, 163828, 98169, 6931,
146           19488, 175760])
147 first_captions = [decoded_captions[elem] for elem
in mask]
148
149 img_idxs = data['val_image_idxs'][mask]          #
the images the captions refer to
150 first_images = [image_from_url(data['val_urls'][j]) for j in img_idxs]
151 #%%
152 for i, (caption, image) in enumerate(zip(
first_captions, first_images)):
153     plt.imshow(image)
154     plt.axis('off')
155     caption_str = caption
156     plt.title(caption_str)
157     plt.show()
158 #%% md
159 # Running the CLIP Model
160
161 First we'll use the pretrained CLIP model to
extract features from the texts and images
separately.
162 #%%
163 device = "mps" if torch.backends.mps.is_available
() else "cpu"
164 clip_model, clip_preprocess = clip.load("ViT-B/32",
device=device)
165 #%%
166 # You can check the model layers by printing the
model.
167 # CLIP's model code is available at https://github.com/openai/CLIP/tree/main/clip

```

```
168 print(clip_model)
169 #%%
170 # First, we encode the captions into vectors in
171 # the shared embedding space.
172 # Since we're using a Transformer as the text
173 # encoder, we need to tokenize the text first.
174 text_tokens = clip.tokenize(first_captions).to(
175     device)
176 with torch.no_grad():
177     text_features = clip_model.encode_text(
178         text_tokens)
179
180 # Sanity check, print the shape
181 print(text_features.shape)
182 #%%
183 # Then, we encode the images into the same
184 # embedding space.
185 processed_images = [
186     clip_preprocess(Image.fromarray(img)).
187     unsqueeze(0)
188     for img in first_images
189 ]
190 images_tensor = torch.cat(processed_images, dim=0
191 ).to(device)
192
193 with torch.no_grad():
194     image_features = clip_model.encode_image(
195         images_tensor)
196
197 # sanity check, print the shape
198 print(image_features.shape)
199 #%% md
200 Open `dl/clip_dino.py` and implement `
201     get_similarity_no_loop` to compute similarity
202     scores between text features and image features.
203     Test your implementation below, you should see
204     relative errors less than 1e-5.
205 #%%
206 from dl.clip_dino import get_similarity_no_loop
207 torch.manual_seed(231)
208 np.random.seed(231)
```

```

197 M, N, D = 5, 6, 10
198
199 test_text_features = torch.randn(N, D)
200 test_image_features = torch.randn(M, D)
201 out = get_similarity_no_loop(test_text_features,
    test_image_features)
202
203 expected_out = np.array([
204     [ 0.1867811 , -0.23494351,  0.44155994, -0.
205     18950461,  0.00100103],
206     [ 0.17905031, -0.25469488, -0.64330417,  0.
207     25097957, -0.09327742],
208     [-0.4407011 , -0.4365381 ,  0.32857686, -0.
209     3765278 ,  0.01049389],
210     [ 0.24815483,  0.42157224, -0.08459304,  0.
211     14132318, -0.26935193],
212     [ 0.02309848, -0.01441101,  0.5469337 ,  0.
213     6018773 ,  0.21581158],
214     [ 0.41579214, -0.014449 , -0.7242257 ,  0.
215     39348006,  0.0822239 ],
216 ]).astype(np.float32)
217
218 print("relative error: ", rel_error(out.numpy(),
219                                     expected_out))
220 #%%
221 # Let's visualize the similarities between our
222 # batch of images and their captions.
223
224 similarities = get_similarity_no_loop(
225     text_features, image_features).cpu().detach().
226     numpy()
227
228 plt.figure(figsize=(20, 14))
229 plt.imshow(similarities, vmin=0.1, vmax=0.3)
230 plt.yticks(range(len(text_features)),
231             first_captions, fontsize=18)
232 plt.xticks([])
233 for i, image in enumerate(first_images):
234     plt.imshow(image, extent=(i - 0.5, i + 0.5, -1
235     .6, -0.6), origin="lower")
236 for x in range(similarities.shape[1]):
```

```
225     for y in range(similarities.shape[0]):  
226         plt.text(x, y, f"{similarities[y, x]:.2f}"  
227             , ha="center", va="center", size=12)  
228  
229 for side in ["left", "top", "right", "bottom"]:  
230     plt.gca().spines[side].set_visible(False)  
231  
232 plt.xlim([-0.5, len(image_features) - 0.5])  
233 plt.ylim([len(text_features) + 0.5, -2])  
234  
235 plt.title("Cosine similarity between text and  
236 image features", size=20)  
237 plt.show()  
238

---

  
239 You will be able to see a high similarity between  
matching image-caption pairs above. We can  
leverage this property to design an image  
classifier that doesn't require any labeled data (i.e., a zero-shot classifier). Each class can be  
represented using an appropriate natural language  
description, and any input image will be  
classified into the class whose description has  
the highest similarity with the image in CLIP's  
embedding space.  
240

---

  
241 #%% md  
242 Implement `clip_zero_shot_classifier` in `dl/  
243 clip_dino.py` and test it below. You should be  
able to see the following predictions:  
244  
245 ['a person', 'an animal', 'an animal', 'food', 'a  
246 person', 'a landscape', 'other', 'other', 'other  
247 ', 'a person']  
248

---

  
249 from dl.clip_dino import clip_zero_shot_classifier  
250  
251 classes = ["a person", "an animal", "food", "a  
252 landscape", "other"]  
253  
254 pred_classes = clip_zero_shot_classifier(  
255
```

```
250     clip_model, clip_preprocess, first_images,
251     classes, device)
252 print(pred_classes)
253 #%% md
254 Run the cell below to visualize the predictions.
255 As you can see, CLIP offers a straightforward way
256 to perform reasonable zero-shot classification
257 across any class taxonomy.
258 #%% md
259 # Visualize the zero shot predictions
260 for i, (pred_class, image) in enumerate(zip(
261     pred_classes, first_images)):
262     plt.imshow(image)
263     plt.axis('off')
264     plt.title(pred_class)
265     plt.show()
266 #%% md
267 # Image Retrieval using CLIP
268 Just as we used CLIP to retrieve the matching
269 class name for each image, we can also use it to
270 retrieve matching images from text inputs (
271 semantic image retrieval). Implement the `CLIPImageRetriever` in `dl/clip_dino.py` and test
272 it by running the two cells below. The expected
273 top 2 outputs for each query are provided in the
274 comments.
275 #%% md
276 from dl.clip_dino import CLIPImageRetriever
277 clip_retriever = CLIPImageRetriever(clip_model,
278     clip_preprocess, first_images, device)
279 #%%
280 query = "sports" # tennis, skateboard
```

```
274 # query = "black and white" # bathroom, zerbias
275 img_indices = clip_retriever.retrieve(query)
276
277 for img_index in img_indices:
278     plt.imshow(first_images[img_index])
279     plt.axis('off')
280     plt.show()
281 %% md
282 **Inline Question 2** -
283
284 CLIP learns to align image and text
representations in a shared latent space using a
contrastive loss. How would you extend this idea
to more than two modalities?
285
286 $\color{blue}{\textit{Your Answer:}}$%
287
288 CLIP's contrastive learning framework can be
extended to more than two modalities by
introducing a dedicated encoder for each modality
(e.g., image, text, audio, video) and projecting
all modalities into a shared latent embedding
space. During training, contrastive losses are
applied between matching samples across modalities
, encouraging representations of the same
underlying content to be close while pushing apart
mismatched samples within a batch. This can be
implemented either by computing pairwise
contrastive losses between all modality pairs or
by anchoring all modalities to a central modality
such as text. Multi-positive contrastive learning
allows all representations belonging to the same
sample to act as positives simultaneously,
improving alignment consistency. As a result, the
model learns modality-agnostic semantic
representations that enable zero-shot
classification and cross-modal retrieval between
any pair of modalities at inference time.
289
290 %% md
291 # DINO
```

292

293 As mentioned earlier, models trained with vanilla contrastive learning methods such as SimCLR and CLIP require very large batch sizes. This makes them computationally expensive and limits their accessibility. Subsequent works, like [BYOL](<https://arxiv.org/abs/2006.07733>), propose an alternative approach that avoids the need for numerous negative samples by using a student-teacher framework. This method performs surprisingly well and was later adopted by [DINO](<https://arxiv.org/abs/2104.14294>) .

294

295 Similar to SimCLR, DINO is trained to maximize the agreement between two vectors derived from different views of the same image. However, unlike SimCLR, DINO uses two separate encoders which are trained differently. The student network is updated via backpropagation to match the outputs of the teacher network. The teacher network is not updated via backpropagation; instead, its weights are updated using an exponential moving average (EMA) of the student's weights. This means that the teacher model evolves more slowly and provides a stable target for the student to learn from.

296

297 Run the cell below to visualize the DINO training pipeline.

298 #%%

299 `from IPython.display import Image as ColabImage`300 `ColabImage(f'dino.gif')`

301 #%%

302 *# first let's get rid of the CLIP model that's currently using memory*303 `del clip_model`304 *# Uncomment the following if you are using GPU runtime*305 `# torch.cuda.empty_cache()`306 `# torch.cuda.ipc_collect()`

307 #%%

308 *# Load smallest dino model. ViT-S/8. Here ViT-S*

```
308 has ~22M parameters and
309 # works on 8x8 patches.
310 dino_model = torch.hub.load('facebookresearch/dino
:main', 'dino_vits8')
311 dino_model.eval().to(device)
312
313 #%%
314 # the image we will be playing around with
315 sample_image = Image.fromarray(first_images[0]).convert("RGB")
316 sample_image
317 #%% md
318 # DINO Attention Maps
319
320 Since the loaded DINO checkpoint is based on the ViT architecture, we can visualize what each attention head is focusing on. The code below generates heatmaps showing which patches of the original image the [CLS] token attends to across the various heads in the final layer. Although this model was trained using a self-supervised objective without any explicit instruction to recognize "structure" in images, still...
321
322 Do you notice any patterns?
323 #%%
324 # Preprocess
325 from torchvision import transforms as T
326 transform = T.Compose([
327     T.Resize((480, 480)),
328     T.ToTensor(),
329     T.Normalize((0.485, 0.456, 0.406), (0.229, 0.
224, 0.225)),
330 ])
331 img_tensor = transform(sample_image)
332 w, h = img_tensor.shape[1:]
333 img_tensor = img_tensor[None].to(device)
334
335 # Extract attention
336 with torch.no_grad():
337     attn = dino_model.get_last_selfattention(
```

```

337 img_tensor)[0, :, 0, 1:]
338 nh, tokens = attn.shape
339 w_feat, h_feat = w // 8, h // 8
340 attn = attn.reshape(nh, w_feat, h_feat)
341 attn = torch.nn.functional.interpolate(attn.
    unsqueeze(0), scale_factor=8, mode="nearest")[0].
    cpu().numpy()
342
343 # Plot attention heads
344 fig, axes = plt.subplots(1, nh, figsize=(3 * nh, 3
    ))
345 for i in range(nh):
346     ax = axes[i] if nh > 1 else axes
347     ax.imshow(attn[i], cmap='inferno')
348     ax.axis('off')
349 plt.show()
350 #%%
351 # Extract patch token features and discard [CLS]
352 # token.
352 with torch.no_grad():
353     all_tokens = dino_model.
354         get_intermediate_layers(img_tensor, n=1)[0] # (1
355         , 1+N, D)
354     patch_tokens = all_tokens[:, 1:, :] # (N, D)
355
356 print(img_tensor.shape)
357 print(all_tokens.shape)
358 print(patch_tokens.shape)
359 #%% md
360
361 **Inline Question 3**
362
363 How do we get the tensor shapes printed above?
363 Explain your answer.
364
365
366 $\color{blue}{\textit Your Answer:}$
367
368 - [1, 3, 480, 480] - This represents a single RGB
368 image of size 480×480.
369 The first dimension (1) is the batch size (one

```

```
369 image), 3 corresponds to the RGB color channels,  
and 480×480 is the spatial resolution.  
370  
371 - [1, 3601, 384]. - DINO uses a Vision Transformer  
with a patch size of 8×8 pixels.  
372 Dividing the image into patches gives: => Image  
width: 480 / 8 = 60, Image height: 480 / 8 = 60  
=> 60x60 = 3600  
373 add the CLS token and get 3601,  
Each token is embedded into a 384-dimensional  
feature vector, where 384 is the fixed embedding  
dimension defined by the DINO ViT architecture. (you  
can find it in the loading dimo model from the  
cell above)  
374  
375 - [1, 3600, 384] - it same as before, we just  
remove the cls token  
376  
377 in summary, the DINO ViT model splits a 480×480  
image into non-overlapping 8×8 patches, producing  
60×60 = 3600 patch tokens. A learnable [CLS] token  
is prepended, yielding 3601 tokens per image.  
Each token is embedded into a 384-dimensional  
feature vector, resulting in tensors of shape (1,  
3601, 384). Removing the CLS token leaves 3600  
patch embeddings of shape (1, 3600, 384), which  
can be reshaped spatially to recover patch-level  
structure.  
378  
379  
380  
381  
382  
383

---

  
384 #%% md  
385 # DINO Features  
386  
387 To understand what the model is encoding in each  
patch, we can visualize the contents of each patch  
token. Since these embeddings are high-  
dimensional and difficult to interpret directly,
```

```
387 we'll use PCA to identify the directions of  
highest variance in the feature space.  
388  
389 In the next cell, we visualize the three principal  
directions of variance in the feature space. This  
reveals the dominant structure that the patch  
embeddings are capturing.  
390 #%%  
391 from sklearn.decomposition import PCA  
392  
393 np.random.seed(231)  
394  
395 # PCA  
396 pca = PCA(n_components=3)  
397 patch_pca = pca.fit_transform(patch_tokens.cpu().  
numpy()[0])  
398  
399 # Normalize PCA components to [0, 1] for RGB  
display  
400 patch_rgb = (patch_pca - patch_pca.min(0)) / (  
patch_pca.max(0) - patch_pca.min(0))  
401  
402 # Reshape to image grid (60x60, 3)  
403 patch_rgb_img = patch_rgb.reshape(60, 60, 3)  
404  
405 # Show as image  
406 plt.figure(figsize=(6, 6))  
407 plt.imshow(patch_rgb_img)  
408 plt.axis('off')  
409 plt.title("Patch Embeddings (PCA → RGB)")  
410 plt.show()  
411 #%% md  
412 **Inline Question 4** -  
413  
414 What kind of structure do you see in the  
visualization above? What does it imply when a  
region consistently appears in a specific color?  
What does it mean when two regions have distinctly  
different color? Remember that PCA reveals the  
directions of highest variance in the feature  
space across all patches. A patch's color reflects
```

```
414 its distinct feature content.  
415  
416  
417 $\\color{blue}{\\textit Your Answer:}$$  
418  
419 The visualization shows clear spatial structure,  
    where patches belonging to the same semantic  
    regions (e.g., the person, the background, the  
    tennis net) form coherent areas with similar  
    colors. This indicates that DINO patch embeddings  
    are locally consistent and capture meaningful  
    visual properties such as object boundaries,  
    textures, and materials, even without supervision.  
420  
421 When a region consistently appears in a specific  
    color, it implies that the patches in that region  
    have similar feature representations and thus  
    project similarly onto the top principal  
    components. In other words, the model encodes  
    those patches as belonging to the same semantic or  
    visual concept. Conversely, when two regions have  
    distinctly different colors, it means their patch  
    embeddings differ significantly in the high-  
    variance directions captured by PCA, indicating  
    different semantic content or visual structure (e.  
    g., foreground vs. background, object vs. texture  
    ). Overall, this demonstrates that DINO features  
    organize image patches in a way that aligns well  
    with object-level and semantic structure.  
422  
423  
424

---

  
425 #%% md  
426 # A Simple Segmentation Model over DINO Features  
427  
428 In the previous section, we saw that DINO features  
    can provide surprisingly good segmentation cues.  
    Now, let's put that idea to the test by training a  
    simple segmentation model on the [DAVIS dataset]  
    (https://davischallenge.org). The DAVIS dataset (Densely Annotated VVideo Segmentation) was created
```

```
428 for video object segmentation tasks. It provides
429 frame-by-frame, pixel-level annotations of objects
430 within videos. For this experiment, we'll train
431 our model using the annotations from just a single
432 frame of a video and see how well it performs on
433 the remaining frames of the same.
434
435 Our model will be intentionally minimal: we'll
436 extract DINO features per patch and train a
437 lightweight per-patch classifier using only the
438 patches from that one annotated frame. Typically,
439 you would train on the full dataset and evaluate
440 on a separate validation set containing different
441 videos. But here, we will test the one-shot
442 capabilities of DINO features.
443
444
445
446
447
448
449
450
451
452
```

```
#%%
```

```
from dl.clip_dino import DavisDataset
```

```
# A helper class to work with DAVIS dataset.
```

```
# It may take ~5 minutes on the first run of this
# cell to download the dataset.
```

```
davis_ds = DavisDataset()
```

```
# Get a specific test video. Do NOT change this
# for submission.
```

```
frames, masks = davis_ds.get_sample(7)
```

```
num_classes = masks.max() + 1
```

```
print(frames.shape, masks.shape, num_classes)
```

```
#%%
```

```
# Get DINO patch features and corresponding class
# labels for a middle frame
```

```
train_fi = 40
```

```
X_train = davis_ds.process_frames(frames[train_fi:
train_fi+1], dino_model, device)[0]
```

```
453 Y_train = davis_ds.process_masks(masks[train_fi:  
    train_fi+1], device)[0]  
454 print(X_train.shape, Y_train.shape)  
455 #%% md  
456 Complete the implementation of the `  
    DINOSegmentation` class in `dl/clip_dino.py`, and  
    test it by running the two cells below. You should  
    achieve a mean IoU greater than 0.45 on the first  
    test frame and greater than 0.50 on the last test  
    frame. To prevent overfitting on the training  
    patch features, consider designing a very  
    lightweight model (e.g., a linear layer or a 2-  
    layer MLP) and applying appropriate weight decay.  
457  
458 You may use GPU runtime to speed up training and  
    evaluation. Make sure to rerun the entire notebook  
    if you change runtime type.  
459 #%%  
460 from dl.clip_dino import DINOSegmentation,  
    compute_iou  
461 torch.manual_seed(231)  
462 np.random.seed(231)  
463 dino_segmentation = DINOSegmentation(device,  
    num_classes)  
464 dino_segmentation.train(X_train, Y_train,  
    num_iters=500)  
465  
466  
467 # Test on first, middle, and last frame  
468 ious = []  
469 test_fis = [0, train_fi, 98]  
470 gt_viz = []  
471 pred_viz = []  
472 for fi in test_fis:  
473     X_test = davis_ds.process_frames(frames[fi:fi+1]  
        , dino_model, device)[0]  
474     Y_test = davis_ds.process_masks(masks[fi:fi+1],  
        device)[0]  
475     Y_pred = dino_segmentation.inference(X_test)  
476     iou = compute_iou(Y_pred, Y_test, num_classes)  
477     ious.append(iou)
```

```

478
479     gt_viz.append(davis_ds.mask_frame_overlay(Y_test
        , frames[fi]))
480     pred_viz.append(davis_ds.mask_frame_overlay(
        Y_pred, frames[fi]))
481
482 gt_viz = np.concatenate(gt_viz, 1)
483 pred_viz = np.concatenate(pred_viz, 1)
484 #%%
485 print(f"Mean IoU on first test frames: {ious[0]:.3f}") # should be >0.45
486 print(f"Mean IoU on last test frames: {ious[2]:.3f}") # should be >0.50
487 #%% md
488 Now let's visualize the results. Run the two cells
    below to display the ground truth and predicted
    segmentation masks for the first, middle, and last
    frames. Note that the middle frame is part of the
    training set, while the other frames are unseen.
489 #%%
490 Image.fromarray(gt_viz)
491 #%%
492 Image.fromarray(pred_viz)
493 #%% md
494 Now run the following three cells to evaluate and
    visualize the entire video. You should achieve a
    mean IoU greater than 0.55. The saved
    visualization video may take some time to process
    in Google Drive, but you can download it to your
    computer and view it locally.
495
496
497 #%%
498 # Run on all frames
499 ious = []
500 gt_viz = []
501 pred_viz = []
502 for fi in range(len(frames)):
503     if fi % 20 == 0:
504         print(f"{fi} / {len(frames)}")
505     X_test = davis_ds.process_frames(frames[fi:fi+1]

```

```

505 ], dino_model, device)[0]
506 Y_test = davis_ds.process_masks(masks[fi:fi+1],
507 device)[0]
508 Y_pred = dino_segmentation.inference(X_test)
509 iou = compute_iou(Y_pred, Y_test, num_classes)
510 ious.append(iou)
511
512 gt_viz.append(davis_ds.mask_frame_overlay(Y_test,
513 , frames[fi]))
514 pred_viz.append(davis_ds.mask_frame_overlay(
515 Y_pred, frames[fi]))
516
517 gt_viz = np.stack(gt_viz) # T x H x W x 3
518 pred_viz = np.stack(pred_viz) # T x H x W x 3
519 final_viz = np.concatenate([gt_viz, pred_viz], -2
520 ) # T x H x 2W x 3
521 #%%
522 print(f"Mean IoU on all frames: {sum(ious) / len(
523 ious):.3f}") # should be >0.55
524
525 #%%
526 import cv2
527
528 def write_video_from_array(array, output_path, fps
529 = 12):
530     T, H, W, _ = array.shape
531     fourcc = cv2.VideoWriter_fourcc(*'mp4v')
532     out = cv2.VideoWriter(output_path, fourcc, fps
533 , (W, H))
534     for i in range(T):
535         frame = array[i]
536         out.write(frame)
537     out.release()
538     print(f"Video saved to {output_path}")
539
540
541 # It might take a while to process in google drive
542 # but you can just download it and watch on your
543 # computer
544 write_video_from_array(final_viz, f"dino_res.mp4")
545 #%% md

```

537 **Inline Question 5** -
538
539 If you train a segmentation model on CLIP ViT's patch features, do you expect it to perform better or worse than DINO? Why should that be the case?
540
541
542
543 \$\color{blue}{\textit{Your Answer:}}\$
544
545 You should expect a segmentation model trained on CLIP ViT patch features to perform worse than one trained on DINO features, especially in this one-shot, per-video segmentation setting.
546
547 Why: DINO is trained with a self-supervised objective that explicitly encourages patch-level consistency and spatial grouping. Its ViT learns to make patches belonging to the same object or region have similar representations, which naturally aligns with segmentation tasks. As a result, DINO patch features tend to be locally coherent, respect object boundaries, and transfer well to dense prediction tasks like segmentation—even with very few labeled examples.
548
549 In contrast, CLIP is trained with a cross-modal contrastive objective at the image-text level. Its primary goal is to align global image representations with text, not to enforce fine-grained spatial structure. While CLIP patch features are semantically meaningful, they are less optimized for spatial coherence, and patch-level representations may vary within the same object as long as the global image-text alignment is preserved. This makes CLIP excellent for zero-shot classification and retrieval, but less effective than DINO for one-shot or few-shot segmentation, where strong patch-level and spatial consistency is crucial.
550

```
551 _____
552 #%%
553 # Submit Your Work
554 # After completing ALL notebooks for this
554 # assignment, run this cell to create submission
554 # files.
555 #
556 # IMPORTANT: Before running this cell:
557 # 1. Make sure all cells in ALL notebooks have
557 # been executed
558 # 2. Manually SAVE all *.ipynb and *.py files
559 #
560 # This will create:
561 # - a2_code_submission.zip (your code)
562 # - a2_inline_submission.pdf (notebook outputs)
563
564 from dl.submit import make_assignment2_submission
565 FOLDERNAME = "" # or keep it as-is if you already
565 # define it earlier
566 make_assignment2_submission("./" + FOLDERNAME )
567
```