Transformer

1 Define the model

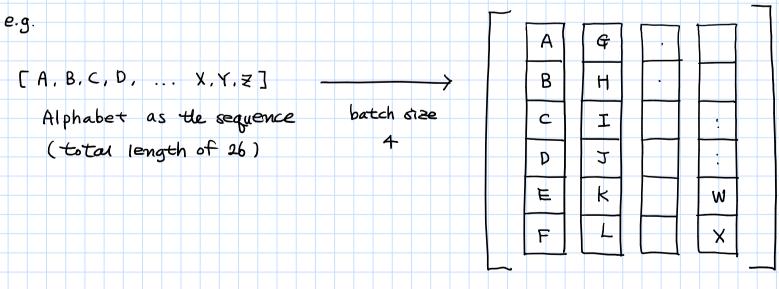
A Language modeling task

- : to assign probability for the likelihood of a given word (or a sequence of words) to follow a sequence of words.
- Sequence of tokens ====> embedding layer passed to
- positional encoding layer
- -Transformer Encoder
 - self-attention layer (square attention mask required)

_ Linear layer

- log-Softmax function

- @ Load & batch data
 - batchify () function
 - : arranges the dataset into columns,
 trimming off any tokens remaining
 after the data has been divided into batches of size batch_size



4 sequences of length 6

```
import torch
import torch. nn as nn
import torch. nn. functional as F
from torch. nn import Transformer Encoder, Transformer Encoder Layer
class Transformer Model (nn. Module):
   def __init__ (self, ntoken: int, d_model: int, nhead: int, d_nid: int,
                  nlayers: int, dropout: float = 0.5):
      super(). \__init__()
       self. model - type = 'Transformer'
      self. pos-encoding = Positional Encoding (d-model, dropout)
      encoder_layers = Transformer Encoder Layer (d_model, nhead, d-hid, dropout)
      self. encoder = nn. Embedding (ntoken, d-model)
       self. d_model = d_model
       self. decoder = nn. Linear (d-model, ntoken)
       self. init _ weights ()
   def init _ weights (self) → None:
       initrange = 0.1
      self. encoder. weight. data. uniform _ (-initrange, initrange)
       self. decoder. bias. data. zero_()
      self. decoder weight data uniform _ (-initrange, initrange)
   def generate_square_subsequent_mask (sz:int) -> Tensor:
       return torch. Triu (torch. ones (SZ.SZ) * float ('-inf'), diagonal = 1)
   def forward (self, src, src_mask):
      src = self. encoder (src) * math. sqrt (self. ninp)
      src = self. pos_encoder (src)
      output = self. transformer _ encoder (src, src_mask)
      output = self. decoder (output)
      return output
```

import math

```
Positional Encoding (nn. Module):
Class
   def __init __ (self, d-model, dropout = 0.1, max_len = 5000):
      super (Positional Encoding, self). __init__()
      self. dropout = nn. Dropout (p = dropout)
      pe = torch. zeros (max_len, 1, d_model)
      position = torch. arange (0, max_len, otype = torch. float). unsqueeze (1)
      div_term = torch. exp (torch. arange (0, d_model, 2). float() *
                            (-math.log(10000.0) / d_model))
      pe[:,0,0:2] = torch. sin (position * div_term)
      pe[:,0,1::2] = torch. Cos (position * div_term)
      pe = pe. unsqueeze (0). transpose (0,1)
      self. register_buffer ('pe', pe)
  def forward (self, x):
      X = X + seif. pe [: X.size(0)]
      return self. dropout (x)
```

```
Load & batch data
 import torch
from torchtext. datasets import WikiText 2
from torchtext. data . utils import get_tokenizer
 from torchtext. Vocab import build _ Vocab _ from _ iterator
 train_iter = WikiText 2 (split = 'train')
 tokenizer = get - tokenizer ('basic _ english')
 vocab = build _ vocab _ from _ iterator (map (tokenizer, train _iter), specials = ["<unk>"])
 Vocab. Set _ default _ index ( Vocab [" <unk>"])
 def data_process (raw_text_iter):
     data = [torch.tensor(vocab(tokenizer(item)), dtype = torch.long) for item in raw_text_iter]
     return torch. Cat (tuple (filter (lambda t: t. nume () > 0, data)))
 train_iter, val_iter, test_iter = WikiText2()
 train - data = data - process (train - iter)
 val_data = data_process (val_iter)
 test_data = data_process (test_iter)
 device = torch device ("cuda" if torch cuda is _ available () else "cpu")
 def batchify (data, bsx):
    # Divide the dataset into bs parts
    seq_len = data size (0) // bsz
    # Trim off any extra elements that wouldn't cleanly fit (remainders)
     data = data [: seq - len * bs=]
     # Evenly divide the data across the bsz batches
     data = data view (bsz, seq len), t(), contiguous ()
     return data to (device)
batch _ size = 20
eval_batch_size = 10
train_data = batchify (train_data, batch_size)
val_data = batchify (val_data, eval_batch_size)
test_data = batchify (test_data, eval_batch_size)
```

```
get_batch()
```

```
- generates a pair of input-target sequences for the transformer model.
```

- it subdivides the source data into chunks of length bptt.

```
def get_batch (source, i):

seq_len = min (bptt, len (source) -1 - i)

data = source [i:i+seq_len]

target = source [i+1:i+1+seq_len]. reshape(-1)

return data, target
```

Initiate an instance

```
ntokens = len(vocab) # the size of Vocabulary

emsize = 200 # embedding dimension

d_hid = 200 # the dimension of feedforward network model in nn. Transformer Encoder

nlayers = 2 # the number of nn. Transformer Encoder Layer in nn. Transformer Encoder

nhead = 2 # the number of heads in the multihead attention models

dropout = 0.2 # the dropout value
```

model = Transformer Model (ntokens, emsize, nhead, n-hid, nlayers, dropout). to (device)

```
Run the model
import time
criterion = nn. Cross Entropy Loss ()
1r = 5.0
optimizer = torch optim. SGD (model. parameters (), Ir = Ir)
scheduler = torch. optim. Ir - scheduler. StepLR (optimizer, 1.0, gamma = 0.95)
def train (model: nn. Module) → None:
    model train ()
    total_1055 = 0.
    log - interval = 200
    start_time = time time ()
    Src_mask = model. generate_square_subsequent_mask (bptt). to (device)
    num_batches = len (train_data) // bptt
    for batch, i in enumerate (range (0, train_data size (0) - 1, bptt)):
        data, targets = get_batch (train_data, i)
        batch _ size = data . size (0)
        if batch_size != bptt:
            src_mask = src_mask [: batch_size, : batch_size]
        output = model (data, src_mask)
        loss = Criterion (output. view (-1, ntokens), targets)
        optimizer. zero - grad()
        loss backward ()
        torch. nn. utils . clip-grad_norm_ (model parameters (), 0.5)
        optimizer. step()
        total_loss += loss.item()
        if batch % log-interval == 0 and batch > 0:
            Ir = scheduler. get_last_lr()[0]
           MS_per_batch = (time.time() - start.time) * 1000 / log_interval
           cur_loss = total_loss / log_interval
            pp1 = math exp (cur_loss)
           print (f' | epoch fepoch: 3d } | f batch: 5d} / f num - batches: 5d} batches!
                  f' |r f |r: 02.2ft | ms/batch | ms_per_batch: 5.2ft |'
                  f'loss { cur_loss: 5.2+ } | pp| {pp|: 8.2+}')
           total_loss = 0
            start_time = time. time()
```

```
def evaluate (model: nn. Module, eval_data: Tensor) -> float:
    model eval () # turn on evaluation mode.
    total_loss = 0.
    Src_mask = generate_square_subsequent_mask (bptt). to (device)
     with torch. no-grad ():
        for i in range (0, eval_data size (0) - 1, bptt):
             data, targets = get _ batch (eval_data, i)
             batch_size = data.size(0)
             if batch _ size != bptt:
                 src_mask = src_mask [: batch_size, : batch_size]
             output = model (data, src_mask)
             output_flat = output. View (-1, ntokens)
             total_loss += batch_size * criterion (output_flat, targets). item()
     return total-loss / (len (eval-data) - 1)
Loop over epochs
 best _ val _ loss = float ('inf')
 epoch = 3
 best _ model = None
for epoch in range (1, epochs +1):
     epoch_start_time = time.time()
    train (model)
     val_loss = evaluate (model, val_data)
    val_ppl = math. exp (val_loss)
    elapsed = time. time() - epoch - start - time
     print ('-' * 89)
     print (f' | end of epoch | epoch : 3d } | time : | elapsed : 5.2f is |
            f'valid loss | val_loss: 5.2ft | valid ppl { val_ppl: 8.2ft')
     print ('-' * 89)
    if val_loss < best _val_loss:
         best_val_loss = val_loss
         best - model = copy. deep copy (model)
    scheduler. Step ()
```

```
Evaluate the best model on the test dataset

test_loss = evaluate (best_model, test_data)

test_ppl = math.exp (test_loss)

print ('=' * 89)

print (f' | End of training | test loss | test_loss: 5.2ft | '

f'test ppl | test_ppl: 8.2fl')

print ('=' * 89)
```