

# Homework 5 - Sequence-to-sequence

若有任何問題，歡迎來信至助教信箱 [ntu-ml-2021spring-ta@googlegroups.com](mailto:ntu-ml-2021spring-ta@googlegroups.com)

- Link to [tutorial\(\)](#) (TBA).

## Sequence-to-Sequence 介紹

- 大多數常見的 seq2seq model 為 encoder-decoder model，主要由兩個部分組成，分別是 encoder 和 decoder，而這兩個部可以使用 recurrent neural network (RNN)或 transformer 來實作，主要是用來解決輸入和輸出的長度不一樣的情況
- **Encoder** 是將一連串的輸入，如文字、影片、聲音訊號等，編碼為單個向量，這單個向量可以想像為是整個輸入的抽象表示，包含了整個輸入的資訊
- **Decoder** 是將 encoder 輸出的單個向量逐步解碼，一次輸出一個結果，直到將最後目標輸出被產生出來為止，每次輸出會影響下一次的輸出，一般會在開頭加入 "< BOS >" 來表示開始解碼，會在結尾輸出 "< EOS >" 來表示輸出結束

## 作業介紹

- 英文翻譯中文
  - 輸入：一句英文 (e.g. tom is a student.)
  - 輸出：中文翻譯 (e.g. 湯姆是個學生。)
- TODO
  - 訓練一個 RNN 模型達到 Seq2seq 翻譯
  - 訓練一個 Transformer 大幅提升效能
  - 實作 Back-translation 大幅提升效能

## 下載和引入需要的函式庫

Load the necessary library: [https://github.com/tensorflow/tensorflow/blob/master/tensorflow/python/tools/import\\_data.py](#)

```
In [ ]: !pip install 'torch>=1.6.0' editdistance matplotlib sacrebleu sacremoses sentencepiece tqdm wandb  
!pip install --upgrade jupyter ipywidgets
```

```
In [ ]: !git clone https://github.com/pytorch/fairseq.git  
!cd fairseq && git checkout 9alc497  
!pip install --upgrade ./fairseq/
```

```
In [ ]:  
import sys  
import pdb  
import pprint  
import logging  
import os  
import random  
  
import torch  
import torch.nn as nn  
import torch.nn.functional as F  
from torch.utils import data  
import numpy as np  
import tqdm.auto as tqdm  
from pathlib import Path  
from argparse import Namespace  
from fairseq import utils  
  
import matplotlib.pyplot as plt
```

## 設定種子

Loading [MathJax]/jax/output/HTML-CSS/fonts/TeX/fontdata.js

```
In [ ]: seed = 73
random.seed(seed)
torch.manual_seed(seed)
if torch.cuda.is_available():
    torch.cuda.manual_seed(seed)
    torch.cuda.manual_seed_all(seed)
np.random.seed(seed)
torch.backends.cudnn.benchmark = False
torch.backends.cudnn.deterministic = True
```

## 資料集介紹

### 英轉繁雙語資料

- [TED2020](#)
  - 原始資料量: 398,066句
  - 處理後資料: 393,980句

### 測試資料

- 資料量: 4,000句
- 中文部分不公開，提供的檔案為假翻譯，全部都是句點。

## 資料下載

### 安裝megatools (optional)

```
In [ ]: #!apt-get install megatools
```

```
Loading [MathJax]/jax/output/HTML-CSS/fonts/TeX/fontdata.js
```

## 下載檔案並解壓縮

```
In [ ]: data_dir = './DATA/rawdata'
dataset_name = 'ted2020'
urls = (
    "https://onedrive.live.com/download?cid=3E549F3B24B238B4&resid=3E549F3B24B238B4%214989&authkey=AGgQ-DaR8eFS11A",
    "https://onedrive.live.com/download?cid=3E549F3B24B238B4&resid=3E549F3B24B238B4%214987&authkey=AA4qP_azsicwZZM",
    # # If the above links die, use the following instead.
    #     "https://www.csie.ntu.edu.tw/~r09922057/ML2021-hw5/ted2020.tgz",
    #     "https://www.csie.ntu.edu.tw/~r09922057/ML2021-hw5/test.tgz",
    # # If the above links die, use the following instead.
    #     "https://mega.nz/#!vEcTCISJ!3Rw0eHTZWPpdHBTbQEjBDikDEdFPr7fI8WxaXK9yZ9U",
    #     "https://mega.nz/#!zNcnGIOJ!oPJX9AvVVsl1jc0SaK6vxP_1FUNTkEcK2WbxJpvjU5Y",
)
file_names = (
    'ted2020.tgz', # train & dev
    'test.tgz', # test
)
prefix = Path(data_dir).absolute() / dataset_name

prefix.mkdir(parents=True, exist_ok=True)
for u, f in zip(urls, file_names):
    path = prefix/f
    if not path.exists():
        if 'mega' in u:
            !megadl {u} --path {path}
        else:
            !wget {u} -O {path}
    if path.suffix == ".tgz":
        !tar -xvf {path} -C {prefix}
    elif path.suffix == ".zip":
        !unzip -o {path} -d {prefix}
!mv {prefix/' raw.en'} {prefix/' train_dev.raw.en'}
!mv {prefix/' raw.zh'} {prefix/' train_dev.raw.zh'}
!mv {prefix/' test.en'} {prefix/' test.raw.en'}
!mv {prefix/' test.zh'} {prefix/' test.raw.zh'}
```

Loading [MathJax]/jax/output/HTML-CSS/fonts/TeX/fontdata.js

## 設定語言

```
In [ ]: src_lang = 'en'  
tgt_lang = 'zh'  
  
data_prefix = f'{prefix} /train_dev.raw'  
test_prefix = f'{prefix} /test.raw'
```

```
In [ ]: !head {data_prefix+'.'+src_lang} -n 5  
!head {data_prefix+'.'+tgt_lang} -n 5
```

## 檔案前處理

Loading [MathJax]/jax/output/HTML-CSS/fonts/TeX/fontdata.js

In [ ]: import re

```

def strQ2B(ustring):
    """把字串全形轉半形"""
    # 參考來源:https://ithelp.ithome.com.tw/articles/10233122
    ss = []
    for s in ustring:
        rstring = ""
        for uchar in s:
            inside_code = ord(uchar)
            if inside_code == 12288: # 全形空格直接轉換
                inside_code = 32
            elif (inside_code >= 65281 and inside_code <= 65374): # 全形字元（除空格）根據關係轉化
                inside_code -= 65248
            rstring += chr(inside_code)
        ss.append(rstring)
    return ''.join(ss)

def clean_s(s, lang):
    if lang == 'en':
        s = re.sub(r"\([^\)]*\)", "", s) # remove ([text])
        s = s.replace('‘', '') # remove ‘
        s = re.sub('([.,;!?\(\)\"])', r' \1 ', s) # keep punctuation
    elif lang == 'zh':
        s = strQ2B(s) # Q2B
        s = re.sub(r"\([^\)]*\)", "", s) # remove ([text])
        s = s.replace('‘', '’')
        s = s.replace('—', '—')
        s = s.replace('“', '”')
        s = s.replace('”', '”')
        s = s.replace('‘’, ’’')
        s = s.replace('‘’, ’’')
        s = re.sub('([.,;!?\(\)\"])', r' \1 ', s) # keep punctuation
    s = ' '.join(s.strip().split())
    return s

def len_s(s, lang):
    if lang == 'zh':
        return len(s)
    return len(s.split())

```

Loading [MathJax]/jax/output/HTML-CSS/fonts/TeX/fontdata.js

```
def clean_corpus(prefix, 11, 12, ratio=9, max_len=1000, min_len=1):
```

```

if Path(f'{prefix}.clean.{11}').exists() and Path(f'{prefix}.clean.{12}').exists():
    print(f'{prefix}.clean.{11} & {12} exists. skipping clean.')
    return
with open(f'{prefix}.{11}', 'r') as 11_in_f:
    with open(f'{prefix}.{12}', 'r') as 12_in_f:
        with open(f'{prefix}.clean.{11}', 'w') as 11_out_f:
            with open(f'{prefix}.clean.{12}', 'w') as 12_out_f:
                for s1 in 11_in_f:
                    s1 = s1.strip()
                    s2 = 12_in_f.readline().strip()
                    s1 = clean_s(s1, 11)
                    s2 = clean_s(s2, 12)
                    s1_len = len_s(s1, 11)
                    s2_len = len_s(s2, 12)
                    if min_len > 0: # remove short sentence
                        if s1_len < min_len or s2_len < min_len:
                            continue
                    if max_len > 0: # remove long sentence
                        if s1_len > max_len or s2_len > max_len:
                            continue
                    if ratio > 0: # remove by ratio of length
                        if s1_len/s2_len > ratio or s2_len/s1_len > ratio:
                            continue
                    print(s1, file=11_out_f)
                    print(s2, file=12_out_f)

```

In [ ]: clean\_corpus(data\_prefix, src\_lang, tgt\_lang)  
clean\_corpus(test\_prefix, src\_lang, tgt\_lang, ratio=-1, min\_len=-1, max\_len=-1)

In [ ]: !head {data\_prefix+'.clean.'+src\_lang} -n 5  
!head {data\_prefix+'.clean.'+tgt\_lang} -n 5

## 切出 train/valid set

In [ ]: valid\_ratio = 0.01 # 3000~4000句就夠了  
train\_ratio = 1 - valid\_ratio

Loading [MathJax]/jax/output/HTML-CSS/fonts/TeX/fontdata.js

```
In [ ]: if (prefix/f'train.clean.{src_lang}).exists() \
and (prefix/f'train.clean.{tgt_lang}).exists() \
and (prefix/f'valid.clean.{src_lang}).exists() \
and (prefix/f'valid.clean.{tgt_lang}).exists():
    print(f'train/valid splits exists. skipping split.')
else:
    line_num = sum(1 for line in open(f'{data_prefix}.clean.{src_lang}'))
    labels = list(range(line_num))
    random.shuffle(labels)
    for lang in [src_lang, tgt_lang]:
        train_f = open(os.path.join(data_dir, dataset_name, f'train.clean.{lang}'), 'w')
        valid_f = open(os.path.join(data_dir, dataset_name, f'valid.clean.{lang}'), 'w')
        count = 0
        for line in open(f'{data_prefix}.clean.{lang}', 'r'):
            if labels[count]/line_num < train_ratio:
                train_f.write(line)
            else:
                valid_f.write(line)
            count += 1
        train_f.close()
        valid_f.close()
```

## Subword Units

翻譯存在的一大問題是未登錄詞(out of vocabulary), 可以使用 subword units 作為斷詞單位來解決。

- 使用 [sentencepiece](#) 套件
- 用 unigram 或 byte-pair encoding (BPE)

```
In [ ]: import sentencepiece as spm
vocab_size = 8000
if (prefix/f' spm{vocab_size}.model').exists():
    print(f'{prefix}/spm{vocab_size}.model exists. skipping spm_train.')
else:
    spm.SentencePieceTrainer.train(
        input=''.join([f'{prefix}/train.clean.{src_lang}', f'{prefix}/valid.clean.{src_lang}', f'{prefix}/train.clean.{tgt_lang}', f'{prefix}/valid.clean.{tgt_lang}']),
        model_prefix=prefix/f' spm{vocab_size}',
        vocab_size=vocab_size,
        character_coverage=1,
        model_type='unigram', # 'bpe' 也可
        input_sentence_size=1e6,
        shuffle_input_sentence=True,
        normalization_rule_name='nmt_nfkc_cf',
    )
```

```
In [ ]: spm_model = spm.SentencePieceProcessor(model_file=str(prefix/f' spm{vocab_size}.model'))
in_tag = {
    'train': 'train.clean',
    'valid': 'valid.clean',
    'test': 'test.raw.clean',
}
for split in ['train', 'valid', 'test']:
    for lang in [src_lang, tgt_lang]:
        out_path = prefix/f'{split}.{lang}'
        if out_path.exists():
            print(f'{out_path} exists. skipping spm_encode.')
        else:
            with open(prefix/f'{split}.{lang}', 'w') as out_f:
                with open(prefix/f'{in_tag[split]}.{lang}', 'r') as in_f:
                    for line in in_f:
                        line = line.strip()
                        tok = spm_model.encode(line, out_type=str)
                        print(''.join(tok), file=out_f)
```

Loading [MathJax]/jax/output/HTML-CSS/fonts/TeX/fontdata.js

```
In [ ]: !head {data_dir+'/' +dataset_name+'train.'+src_lang} -n 5  
!head {data_dir+'/' +dataset_name+'train.'+tgt_lang} -n 5
```

## 用 fairseq 將資料轉為 binary

```
In [ ]: binpath = Path('./DATA/data-bin', dataset_name)  
if binpath.exists():  
    print(binpath, "exists, will not overwrite!")  
else:  
    !python -m fairseq_cli.preprocess []  
    --source-lang {src_lang}\br/>    --target-lang {tgt_lang}\br/>    --trainpref {prefix/'train'}\br/>    --validpref {prefix/'valid'}\br/>    --testpref {prefix/'test'}\br/>    --destdir {binpath}\br/>    --joined-dictionary\br/>    --workers 2
```

## 實驗的參數設定表

Loading [MathJax]/jax/output/HTML-CSS/fonts/TeX/fontdata.js

```
In [ ]: config = Namespace(  
    datadir = "./DATA/data-bin/ted2020",  
    savedir = "./checkpoints/rnn",  
    source_lang = "en",  
    target_lang = "zh",  
  
    # cpu threads when fetching & processing data.  
    num_workers=2,  
    # batch size in terms of tokens. gradient accumulation increases the effective batchsize.  
    max_tokens=8192,  
    accum_steps=2,  
  
    # the lr is calculated from Noam lr scheduler. you can tune the maximum lr by this factor.  
    lr_factor=2.,  
    lr_warmup=4000,  
  
    # clipping gradient norm helps alleviate gradient exploding  
    clip_norm=1.0,  
  
    # maximum epochs for training  
    max_epoch=30,  
    start_epoch=1,  
  
    # beam size for beam search  
    beam=5,  
    # generate sequences of maximum length ax + b, where x is the source length  
    max_len_a=1.2,  
    max_len_b=10,  
    # when decoding, post process sentence by removing sentencepiece symbols.  
    post_process = "sentencepiece",  
  
    # checkpoints  
    keep_last_epochs=5,  
    resume=None, # if resume from checkpoint name (under config.savedir)  
  
    # logging  
    use_wandb=False,  
)
```

Loading [MathJax]/jax/output/HTML-CSS/fonts/TeX/fontdata.js

# Logging

- logging 套件紀錄一般訊息
- wandb 紀錄續練過程 loss, bleu, model weight 等等

```
In [ ]: logging.basicConfig(  
        format"%(asctime)s | %(levelname)s | %(name)s | %(message)s",  
        datefmt"%Y-%m-%d %H:%M:%S",  
        level="INFO", # "DEBUG" "WARNING" "ERROR"  
        stream=sys.stdout,  
)  
proj = "hw5.seq2seq"  
logger = logging.getLogger(proj)  
if config.use_wandb:  
    import wandb  
    wandb.init(project=proj, name=Path(config.savedir).stem, config=config)
```

# CUDA環境

```
In [ ]: cuda_env = utils.CudaEnvironment()  
utils.CudaEnvironment.pretty_print_cuda_env_list([cuda_env])  
device = torch.device('cuda:0' if torch.cuda.is_available() else 'cpu')
```

# 讀取資料集

Loading [MathJax]/jax/output/HTML-CSS/fonts/TeX/fontdata.js

## 借用 fairseq 的 TranslationTask

- 用來讀進上面 binarized 的檔案
- 有現成的 data iterator (dataloader)
- 字典 task.source\_dictionary 和 task.target\_dictionary 也很好用
- 有實做 beam search

```
In [ ]: from fairseq.tasks.translation import TranslationConfig, TranslationTask
```

```
## setup task
task_cfg = TranslationConfig(
    data=config.datadir,
    source_lang=config.source_lang,
    target_lang=config.target_lang,
    train_subset="train",
    required_seq_len_multiple=8,
    dataset_impl="mmap",
    upsample_primary=1,
)
task = TranslationTask.setup_task(task_cfg)
```

```
In [ ]: logger.info("loading data for epoch 1")
task.load_dataset(split="train", epoch=1, combine=True) # combine if you have back-translation data.
task.load_dataset(split="valid", epoch=1)
```

Loading [MathJax]/jax/output/HTML-CSS/fonts/TeX/fontdata.js

```
In [ ]: sample = task.dataset("valid")[1]
pprint.pprint(sample)
pprint.pprint(
    "Source: " + \
    task.source_dictionary.string(
        sample['source'],
        config.post_process,
    )
)
pprint.pprint(
    "Target: " + \
    task.target_dictionary.string(
        sample['target'],
        config.post_process,
    )
)
```

## Dataset Iterator

- 將每個 batch 控制在 N 個 token 讓 GPU 記憶體更有效被利用
- 讓 training set 每個 epoch 有不同 shuffling
- 濾掉長度太長的句子
- 將每個 batch 內的句子 pad 成一樣長，好讓 GPU 平行運算
- 加上 eos 並 shift 一格
  - teacher forcing: 為了訓練模型根據prefix生成下個字，decoder的輸入會是輸出目標序列往右shift一格。
  - 一般是會在輸入開頭加個bos token (如下圖)
  - fairseq 則是直接把 eos 挪到 beginning，訓練起來效果其實差不多。例如:

```
# 輸出目標 (target) 和 Decoder輸入 (prev_output_tokens):
eos = 2
target = 419, 711, 238, 888, 792, 60, 968, 8, 2
prev_output_tokens = 2, 419, 711, 238, 888, 792, 60, 968, 8
```

Loading [MathJax]/jax/output/HTML-CSS/fonts/TeX/fontdata.js

```
In [ ]: def load_data_iterator(task, split, epoch=1, max_tokens=4000, num_workers=1, cached=True):
    batch_iterator = task.get_batch_iterator(
        dataset=task.dataset(split),
        max_tokens=max_tokens,
        max_sentences=None,
        max_positions=utils.resolve_max_positions(
            task.max_positions(),
            max_tokens,
        ),
        ignore_invalid_inputs=True,
        seed=seed,
        num_workers=num_workers,
        epoch=epoch,
        disable_iterator_cache=not cached,
        # Set this to False to speed up. However, if set to False, changing max_tokens beyond
        # first call of this method has no effect.
    )
    return batch_iterator

demo_epoch_obj = load_data_iterator(task, "valid", epoch=1, max_tokens=20, num_workers=1, cached=False)
demo_iter = demo_epoch_obj.next_epoch_itr(shuffle=True)
sample = next(demo_iter)
sample
```

Loading [MathJax]/jax/output/HTML-CSS/fonts/TeX/fontdata.js

- 每個 batch 是一個字典，key 是字串，value 是 Tensor，內容說明如下

```
batch = {
    "id": id, # 每個 example 的 id
    "nsentences": len(samples), # batch size 句子數
    "ntokens": ntokens, # batch size 字數
    "net_input": {
        "src_tokens": src_tokens, # 來源語言的序列
        "src_lengths": src_lengths, # 每句話沒有 pad 過的長度
        "prev_output_tokens": prev_output_tokens, # 上面提到右 shift 一格後的目標序列
    },
    "target": target, # 目標序列
}
```

## 定義模型架構

- 我們一樣繼承 fairseq 的 encoder, decoder 和 model，這樣測試階段才能直接用他寫好的 beam search 函式

```
In [ ]: from fairseq.models import (
    FairseqEncoder,
    FairseqIncrementalDecoder,
    FairseqEncoderDecoderModel
)
```

## Encoder 編碼器

Loading [MathJax]/jax/output/HTML-CSS/fonts/TeX/fontdata.js

- seq2seq 模型的編碼器為 RNN 或 Transformer Encoder，以下說明以 RNN 為例，Transformer 略有不同。對於每個輸入，Encoder 會輸出一個向量和一個隱藏狀態(hidden state)，並將隱藏狀態用於下一個輸入。換句話說，Encoder 會逐步讀取輸入序列，並在每個 timestep 輸出單個向量，以及在最後 timestep 輸出最終隱藏狀態(content vector)
- 參數：
  - *args*
    - *encoder\_embed\_dim* 是 embedding 的維度，主要將 one-hot vector 的單詞向量壓縮到指定的維度，主要是為了降維和濃縮資訊的功用
    - *encoder\_ffn\_embed\_dim* 是 RNN 輸出和隱藏狀態的維度(hidden dimension)
    - *encoder\_layers* 是 RNN 要疊多少層
    - *dropout* 是決定有多少的機率會將某個節點變為 0，主要是為了防止 overfitting，一般來說是在訓練時使用，測試時則不使用
  - *dictionary*: fairseq 幫我們做好的 dictionary. 在此用來得到 padding index，好用來得到 encoder padding mask.
  - *embed\_tokens*: 事先做好的詞嵌入 (nn.Embedding)
- 輸入：
  - *src\_tokens*: 英文的整數序列 e.g. 1, 28, 29, 205, 2
- 輸出：
  - *outputs*: 最上層 RNN 每個 timestep 的輸出，後續可以用 Attention 再進行處理
  - *final\_hiddens*: 每層最終 timestep 的隱藏狀態，將傳遞到 Decoder 進行解碼
  - *encoder\_padding\_mask*: 告訴我們哪些是位置的資訊不重要。

Loading [MathJax]/jax/output/HTML-CSS/fonts/TeX/fontdata.js

```
In [ ]: class RNNEncoder(FairseqEncoder):
    def __init__(self, args, dictionary, embed_tokens):
        super().__init__(dictionary)
        self.embed_tokens = embed_tokens

        self.embed_dim = args.encoder_embed_dim
        self.hidden_dim = args.encoder_ffn_embed_dim
        self.num_layers = args.encoder_layers

        self.dropout_in_module = nn.Dropout(args.dropout)
        self.rnn = nn.GRU(
            self.embed_dim,
            self.hidden_dim,
            self.num_layers,
            dropout=args.dropout,
            batch_first=False,
            bidirectional=True
        )
        self.dropout_out_module = nn.Dropout(args.dropout)

        self.padding_idx = dictionary.pad()

    def combine_bidir(self, outs, bsz: int):
        out = outs.view(self.num_layers, 2, bsz, -1).transpose(1, 2).contiguous()
        return out.view(self.num_layers, bsz, -1)

    def forward(self, src_tokens, **unused):
        bsz, seqlen = src_tokens.size()

        # get embeddings
        x = self.embed_tokens(src_tokens)
        x = self.dropout_in_module(x)

        # B x T x C -> T x B x C
        x = x.transpose(0, 1)

        # 過雙向RNN
        h0 = x.new_zeros(2 * self.num_layers, bsz, self.hidden_dim)
        x, final_hiddens = self.rnn(x, h0)
        # outputs - [sequence len, batch size, hid dim * directions] 是最上層RNN的輸出

```

Loading [MathJax]/jax/output/HTML-CSS/fonts/TeX/fontdata.js  
# outputs - [sequence len, batch size, hid dim \* directions] 是最上層RNN的輸出

```
# hidden = [num_layers * directions, batch size , hid dim]

# 因為 Encoder 是雙向的RNN，所以需要將同一層兩個方向的 hidden state 接在一起
final_hiddens = self.combine_bidir(final_hiddens, bsz)
# hidden = [num_layers x batch x num_directions*hidden]

encoder_padding_mask = src_tokens.eq(self.padding_idx).t()

return tuple(
(
    outputs, # seq_len x batch x hidden
    final_hiddens, # num_layers x batch x num_directions*hidden
    encoder_padding_mask, # seq_len x batch
)
)

def reorder_encoder_out(self, encoder_out, new_order):
    # 這個beam search時會用到，意義並不是很重要
    return tuple(
(
    encoder_out[0].index_select(1, new_order),
    encoder_out[1].index_select(1, new_order),
    encoder_out[2].index_select(1, new_order),
)
)
```

## Attention

Loading [MathJax]/jax/output/HTML-CSS/fonts/TeX/fontdata.js

- 當輸入過長，或是單獨靠 “content vector” 無法取得整個輸入的意思時，用 Attention Mechanism 來提供 Decoder 更多的資訊
- 根據現在 **Decoder embeddings**，去計算在 **Encoder outputs** 中，那些與其有較高的關係，根據關係的數值來把 Encoder outputs 平均起來作為 **Decoder RNN** 的輸入
- 常見 Attention 的實作是用 Neural Network / Dot Product 來算 **query** (decoder embeddings) 和 **key** (Encoder outputs) 之間的關係，再對所有算出來的數值做 **softmax** 得到分佈，最後根據這個分佈對 **values** (Encoder outputs) 做 **weight sum**
- 參數：
  - *input\_embed\_dim*: key 的維度，應是 decoder 要做 attend 時的向量的維度
  - *source\_embed\_dim*: query 的維度，應是要被 attend 的向量(encoder outputs)的維度
  - *output\_embed\_dim*: value 的維度，應是做完 attention 後，下一層預期的向量維度
- 輸入：
  - *inputs*: 就是 key，要 attend 別人的向量
  - *encoder\_outputs*: 是 query/value，被 attend 的向量
  - *encoder\_padding\_mask*: 告訴我們哪些是位置的資訊不重要。
- 輸出：
  - *output*: 做完 attention 後的 context vector
  - *attention score*: attention 的分布

Loading [MathJax]/jax/output/HTML-CSS/fonts/TeX/fontdata.js

```
In [ ]: class AttentionLayer(nn.Module):
    def __init__(self, input_embed_dim, source_embed_dim, output_embed_dim, bias=False):
        super().__init__()

        self.input_proj = nn.Linear(input_embed_dim, source_embed_dim, bias=bias)
        self.output_proj = nn.Linear(
            input_embed_dim + source_embed_dim, output_embed_dim, bias=bias
        )

    def forward(self, inputs, encoder_outputs, encoder_padding_mask):
        # inputs: T, B, dim
        # encoder_outputs: S x B x dim
        # padding mask: S x B

        # convert all to batch first
        inputs = inputs.transpose(1, 0) # B, T, dim
        encoder_outputs = encoder_outputs.transpose(1, 0) # B, S, dim
        encoder_padding_mask = encoder_padding_mask.transpose(1, 0) # B, S

        # 投影到encoder_outputs的維度
        x = self.input_proj(inputs)

        # 計算attention
        # (B, T, dim) x (B, dim, S) = (B, T, S)
        attn_scores = torch.bmm(x, encoder_outputs.transpose(1, 2))

        # 擋住padding位置的attention
        if encoder_padding_mask is not None:
            # 利用broadcast B, S -> (B, 1, S)
            encoder_padding_mask = encoder_padding_mask.unsqueeze(1)
            attn_scores = (
                attn_scores.float()
                .masked_fill_(encoder_padding_mask, float("-inf"))
                .type_as(attn_scores)
            ) # FP16 support: cast to float and back

        # 在source對應維度softmax
        attn_scores = F.softmax(attn_scores, dim=-1)
```

Loading [MathJax]/jax/output/HTML-CSS/TeX/fontdata.js  
 $\#_形狀(B, T, S) \times (B, S, dim) = (B, T, dim)$  加權平均  
 $x = \text{torch.bmm}(\text{attn\_scores}, \text{encoder\_outputs})$

```

# (B, T, dim)
x = torch.cat((x, inputs), dim=-1)
x = torch.tanh(self.output_proj(x)) # concat + linear + tanh

# 回復形狀 (B, T, dim) -> (T, B, dim)
return x.transpose(1, 0), attn_scores

```

## Decoder 解碼器

- 解碼器的 hidden states 會用編碼器最終隱藏狀態來初始化(content vector)
- 解碼器同時也根據目前 timestep 的輸入(也就是前幾個 timestep 的 output), 改變 hidden states, 並輸出結果
- 如果加入 attention 可以使表現更好
- 我們把 seq2seq 步驟寫在解碼器裡, 好讓等等 Seq2Seq 這個型別可以通用 RNN 和 Transformer, 而不用再改寫
- 參數:
  - args*
    - decoder\_embed\_dim 是解碼器 embedding 的維度, 類同 encoder\_embed\_dim,
    - decoder\_ffn\_embed\_dim 是解碼器 RNN 的隱藏維度, 類同 encoder\_ffn\_embed\_dim
    - decoder\_layers 解碼器 RNN 的層數
    - share\_decoder\_input\_output\_embed 通常 decoder 最後輸出的投影矩陣會和輸入 embedding 共用參數
  - dictionary*: fairseq 幫我們做好的 dictionary.
  - embed\_tokens*: 事先做好的詞嵌入(nn.Embedding)
- 輸入:
  - prev\_output\_tokens*: 英文的整數序列 e.g. 1, 28, 29, 205, 2 已經 shift 一格的 target
  - encoder\_out*: 編碼器的輸出
  - incremental\_state*: 這是測試階段為了加速, 所以會記錄每個 timestep 的 hidden state 詳見 forward
- 輸出:
  - outputs*: decoder 每個 timestep 的 logits, 還沒經過 softmax 的分布
  - extra*: 沒用到

```
In [ ]: class RNNDncoder(FairseqIncrementalDecoder):
    def __init__(self, args, dictionary, embed_tokens):
        super().__init__(dictionary)
        self.embed_tokens = embed_tokens

        assert args.decoder_layers == args.encoder_layers, f"""seq2seq rnn requires that encoder
and decoder have same layers of rnn. got: {args.encoder_layers, args.decoder_layers}"""
        assert args.decoder_ffn_embed_dim == args.encoder_ffn_embed_dim*2, f"""seq2seq-rnn requires
that decoder hidden to be 2*encoder hidden dim. got: {args.decoder_ffn_embed_dim, args.encoder_ffn_embed_dim*2}"""

        self.embed_dim = args.decoder_embed_dim
        self.hidden_dim = args.decoder_ffn_embed_dim
        self.num_layers = args.decoder_layers

        self.dropout_in_module = nn.Dropout(args.dropout)
        self.rnn = nn.GRU(
            self.embed_dim,
            self.hidden_dim,
            self.num_layers,
            dropout=args.dropout,
            batch_first=False,
            bidirectional=False
        )
        self.attention = AttentionLayer(
            self.embed_dim, self.hidden_dim, self.embed_dim, bias=False
        )
        # self.attention = None
        self.dropout_out_module = nn.Dropout(args.dropout)

        if self.hidden_dim != self.embed_dim:
            self.project_out_dim = nn.Linear(self.hidden_dim, self.embed_dim)
        else:
            self.project_out_dim = None

        if args.share_decoder_input_output_embed:
            self.output_projection = nn.Linear(
                self.embed_tokens.weight.shape[1],
                self.embed_tokens.weight.shape[0],
            )
```

Loading [MathJax]/jax/output/HTML-CSS/fonts/TeX/fontdata.js

```

        self.output_projection.weight = self.embed_tokens.weight
    else:
        self.output_projection = nn.Linear(
            self.output_embed_dim, len(dictionary), bias=False
        )
        nn.init.normal_(
            self.output_projection.weight, mean=0, std=self.output_embed_dim ** -0.5
        )

def forward(self, prev_output_tokens, encoder_out, incremental_state=None, **unused):
    # 取出encoder的輸出
    encoder_outputs, encoder_hiddens, encoder_padding_mask = encoder_out
    # outputs:          seq_len x batch x num_directions*hidden
    # encoder_hiddens: num_layers x batch x num_directions*encoder_hidden
    # padding_mask:     seq_len x batch

    if incremental_state is not None and len(incremental_state) > 0:
        # 有上個timestep留下的資訊，讀進來就可以繼續decode，不用從bos重來
        prev_output_tokens = prev_output_tokens[:, -1:]
        cache_state = self.get_incremental_state(incremental_state, "cached_state")
        prev_hiddens = cache_state["prev_hiddens"]
    else:
        # 沒有incremental state代表這是training或者是test time的第一步
        # 準備seq2seq: 把encoder_hiddens pass進去decoder的hidden states
        prev_hiddens = encoder_hiddens

    bsz, seqlen = prev_output_tokens.size()

    # embed tokens
    x = self.embed_tokens(prev_output_tokens)
    x = self.dropout_in_module(x)

    # B x T x C -> T x B x C
    x = x.transpose(0, 1)

    # 做decoder-to-encoder attention
    if self.attention is not None:
        x, attn = self.attention(x, encoder_outputs, encoder_padding_mask)

    # 過單向RNN
    Loading [MathJax]/jax/output/HTML-CSS/fonts/TeX/fontdata.js (x, prev_hiddens)
    # outputs = [sequence len, batch size, hid dim]

```

```

# hidden = [num_layers * directions, batch size , hid dim]
x = self.dropout_out_module(x)

# 投影到embedding size (如果hidden 和embed size不一樣, 然後share_embedding又設成True, 需要額外project一次)
if self.project_out_dim != None:
    x = self.project_out_dim(x)

# 投影到vocab size 的分佈
x = self.output_projection(x)

# T x B x C -> B x T x C
x = x.transpose(1, 0)

# 如果是Incremental, 記錄這個timestep的hidden states, 下個timestep讀回來
cache_state = {
    "prev_hiddens": final_hiddens,
}
self.set_incremental_state(incremental_state, "cached_state", cache_state)

return x, None

def reorder_incremental_state(
    self,
    incremental_state,
    new_order,
):
    # 這個beam search時會用到, 意義並不是很重要
    cache_state = self.get_incremental_state(incremental_state, "cached_state")
    prev_hiddens = cache_state["prev_hiddens"]
    prev_hiddens = [p.index_select(0, new_order) for p in prev_hiddens]
    cache_state = {
        "prev_hiddens": torch.stack(prev_hiddens),
    }
    self.set_incremental_state(incremental_state, "cached_state", cache_state)
    return

```

## Seq2Seq

- 由 **Encoder** 和 **Decoder** 組成
- 接收輸入並傳給 **Encoder**
- 將 **Encoder** 的輸出傳給 **Decoder**
- **Decoder** 根據前幾個 timestep 的輸出和 **Encoder** 輸出進行解碼
- 當解碼完成後，將 **Decoder** 的輸出傳回

```
In [ ]: class Seq2Seq(FairseqEncoderDecoderModel):  
    def __init__(self, args, encoder, decoder):  
        super().__init__(encoder, decoder)  
        self.args = args  
  
    def forward(  
        self,  
        src_tokens,  
        src_lengths,  
        prev_output_tokens,  
        return_all_hiddens: bool = True,  
    ):  
        """  
        Run the forward pass for an encoder-decoder model.  
        """  
        encoder_out = self.encoder(  
            src_tokens, src_lengths=src_lengths, return_all_hiddens=return_all_hiddens  
        )  
        logits, extra = self.decoder(  
            prev_output_tokens,  
            encoder_out=encoder_out,  
            src_lengths=src_lengths,  
            return_all_hiddens=return_all_hiddens,  
        )  
        return logits, extra
```

Loading [MathJax]/jax/output/HTML-CSS/fonts/TeX/fontdata.js

# 模型初始化

Loading [MathJax]/jax/output/HTML-CSS/fonts/TeX/fontdata.js

```
In [ ]: # # HINT: transformer 架構
# from fairseq.models.transformer import (
#     TransformerEncoder,
#     TransformerDecoder,
# )

def build_model(args, task):
    """按照參數設定建置模型"""
    src_dict, tgt_dict = task.source_dictionary, task.target_dictionary

    # 詞嵌入
    encoder_embed_tokens = nn.Embedding(len(src_dict), args.encoder_embed_dim, src_dict.pad())
    decoder_embed_tokens = nn.Embedding(len(tgt_dict), args.decoder_embed_dim, tgt_dict.pad())

    # 編碼器與解碼器
    # TODO: 替換成 TransformerEncoder 和 TransformerDecoder
    encoder = RNNEncoder(args, src_dict, encoder_embed_tokens)
    decoder = RNNDecoder(args, tgt_dict, decoder_embed_tokens)

    # 序列到序列模型
    model = Seq2Seq(args, encoder, decoder)

    # 序列到序列模型的初始化很重要 需要特別處理
    def init_params(module):
        from fairseq.modules import MultiheadAttention
        if isinstance(module, nn.Linear):
            module.weight.data.normal_(mean=0.0, std=0.02)
            if module.bias is not None:
                module.bias.data.zero_()
        if isinstance(module, nn.Embedding):
            module.weight.data.normal_(mean=0.0, std=0.02)
            if module.padding_idx is not None:
                module.weight.data[module.padding_idx].zero_()
        if isinstance(module, MultiheadAttention):
            module.q_proj.weight.data.normal_(mean=0.0, std=0.02)
            module.k_proj.weight.data.normal_(mean=0.0, std=0.02)
            module.v_proj.weight.data.normal_(mean=0.0, std=0.02)
        if isinstance(module, nn.RNNBase):
            for name, param in module.named_parameters():
                if "weight" in name or "bias" in name:
                    param.data.uniform_(-0.1, 0.1)
```

Loading [MathJax]/jax/output/HTML-CSS/fonts/TeX/fontdata.js

```
# 初始化模型  
model.apply(init_params)  
return model
```

## 設定模型相關參數

參考參數

model	embedding dim	encoder ffn	encoder layers	decoder ffn	decoder layers
RNN	256	512	1	1024	1
Transformer	256	1024	4	1024	4

Strong baseline 用的參數可以參考 [Attention is all you need](#) 的 Table 3 的 transformer-base

Loading [MathJax]/jax/output/HTML-CSS/fonts/TeX/fontdata.js

```
In [ ]: arch_args = Namespace(
    encoder_embed_dim=256,
    encoder_ffn_embed_dim=512,
    encoder_layers=1,
    decoder_embed_dim=256,
    decoder_ffn_embed_dim=1024,
    decoder_layers=1,
    share_decoder_input_output_embed=True,
    dropout=0.3,
)

## HINT: 補上Transformer用的參數
# def add_transformer_args(args):
#     args.encoder_attention_heads=4
#     args.encoder_normalize_before=True

#     args.decoder_attention_heads=4
#     args.decoder_normalize_before=True

#     args.activation_fn="relu"
#     args.max_source_positions=1024
#     args.max_target_positions=1024

#     # 補上我們沒有設定的Transformer預設參數
#     from fairseq.models.transformer import base_architecture
#     base_architecture(arch_args)

# add_transformer_args(arch_args)
```

```
In [ ]: if config.use_wandb:
    wandb.config.update(vars(arch_args))
```

```
In [ ]: model = build_model(arch_args, task)
logger.info(model)
```

## Optimization 最佳化

## Loss: Label Smoothing Regularization

- 讓模型學習輸出較不集中的分佈，防止模型過度自信
- 有時候Ground Truth並非唯一答案，所以在算loss時，我們會保留一部份機率給正確答案以外的label
- 可以有效防止過度擬合

code source ([https://fairseq.readthedocs.io/en/latest/\\_modules/fairseq/criterions/label\\_smoothed\\_cross\\_entropy.html](https://fairseq.readthedocs.io/en/latest/_modules/fairseq/criterions/label_smoothed_cross_entropy.html))

```
In [ ]: class LabelSmoothedCrossEntropyCriterion(nn.Module):
    def __init__(self, smoothing, ignore_index=None, reduce=True):
        super().__init__()
        self.smoothing = smoothing
        self.ignore_index = ignore_index
        self.reduce = reduce

    def forward(self, lprobs, target):
        if target.dim() == lprobs.dim() - 1:
            target = target.unsqueeze(-1)
        # nll: Negative log likelihood, 當目標是one-hot時的cross-entropy loss. 以下同 F.nll_loss
        nll_loss = -lprobs.gather(dim=-1, index=target)
        # 將一部分正確答案的機率分配給其他label 所以當計算cross-entropy時等於把所有label的log prob加起來
        smooth_loss = -lprobs.sum(dim=-1, keepdim=True)
        if self.ignore_index is not None:
            pad_mask = target.eq(self.ignore_index)
            nll_loss.masked_fill_(pad_mask, 0.0)
            smooth_loss.masked_fill_(pad_mask, 0.0)
        else:
            nll_loss = nll_loss.squeeze(-1)
            smooth_loss = smooth_loss.squeeze(-1)
        if self.reduce:
            nll_loss = nll_loss.sum()
            smooth_loss = smooth_loss.sum()
        # 計算cross-entropy時 加入分配給其他label的loss
        eps_i = self.smoothing / lprobs.size(-1)
        loss = (1.0 - self.smoothing) * nll_loss + eps_i * smooth_loss
        return loss

# 一般都用0.1效果就很好了
criterion = LabelSmoothedCrossEntropyCriterion(
    smoothing=0.1,
    ignore_index=task.target_dictionary.pad(),
)
```

Loading [MathJax]/jax/output/HTML-CSS/fonts/TeX/fontdata.js

## Optimizer: Adam + lr scheduling

Inverse square root 排程對於訓練 Transformer 時的穩定性很重要，後來也用在 RNN 上。根據底下公式來更新 learning rate，前期線性增長，後期根據更新步數方根的倒數來遞減。

$$lrate = d_{\text{model}}^{-0.5} \cdot \min (step\_num^{-0.5}, step\_num \cdot warmup\_steps^{-1.5})$$

code source (<https://nlp.seas.harvard.edu/2018/04/03/attention.html>)

Loading [MathJax]/jax/output/HTML-CSS/fonts/TeX/fontdata.js

```
In [ ]: class NoamOpt:  
    """Optim wrapper that implements rate.  
    def __init__(self, model_size, factor, warmup, optimizer):  
        self.optimizer = optimizer  
        self._step = 0  
        self.warmup = warmup  
        self.factor = factor  
        self.model_size = model_size  
        self._rate = 0  
  
    @property  
    def param_groups(self):  
        return self.optimizer.param_groups  
  
    def multiply_grads(self, c):  
        """Multiplies grads by a constant *c*.  
        for group in self.param_groups:  
            for p in group['params']:  
                if p.grad is not None:  
                    p.grad.data.mul_(c)  
  
    def step(self):  
        """Update parameters and rate  
        self._step += 1  
        rate = self.rate()  
        for p in self.param_groups:  
            p['lr'] = rate  
        self._rate = rate  
        self.optimizer.step()  
  
    def rate(self, step = None):  
        """Implement `lrate` above  
        if step is None:  
            step = self._step  
        return 0 if not step else self.factor * \  
            (self.model_size ** (-0.5)) *  
            min(step ** (-0.5), step * self.warmup ** (-1.5)))
```

Loading [MathJax]/jax/output/HTML-CSS/fonts/TeX/fontdata.js

## 排程視覺化

```
In [ ]: optimizer = NoamOpt(
    model_size=arch_args.encoder_embed_dim,
    factor=config.lr_factor,
    warmup=config.lr_warmup,
    optimizer=torch.optim.AdamW(model.parameters(), lr=0, betas=(0.9, 0.98), eps=1e-09, weight_decay=0.0001))
plt.plot(np.arange(1, 100000), [optimizer.rate(i) for i in range(1, 100000)])
plt.legend([f'{optimizer.model_size}: {optimizer.warmup}'])
None
```

## 訓練步驟

### Training 訓練

Loading [MathJax]/jax/output/HTML-CSS/fonts/TeX/fontdata.js

```
In [ ]: from fairseq.data import iterators
from torch.cuda.amp import GradScaler, autocast

def train_one_epoch(epoch_itr, model, task, criterion, optimizer, accum_steps=1):
    itr = epoch_itr.next_epoch_itr(shuffle=True)
    itr = iterators.GroupedIterator(itr, accum_steps) # 梯度累積: 每 accum_steps 個 sample 更新一次

    stats = {"loss": []}
    scaler = GradScaler() # 混和精度訓練 automatic mixed precision (amp)

    model.train()
    progress = tqdm.tqdm(itr, desc=f"train epoch {epoch_itr.epoch}", leave=False)
    for samples in progress:
        model.zero_grad()
        accum_loss = 0
        sample_size = 0
        # 梯度累積: 每 accum_steps 個 sample 更新一次
        for i, sample in enumerate(samples):
            if i == 1:
                # emptying the CUDA cache after the first step can reduce the chance of OOM
                torch.cuda.empty_cache()

            sample = utils.move_to_cuda(sample, device=device)
            target = sample["target"]
            sample_size_i = sample["ntokens"]
            sample_size += sample_size_i

            # 混和精度訓練
            with autocast():
                net_output = model.forward(**sample["net_input"])
                lprobs = F.log_softmax(net_output[0], -1)
                loss = criterion(lprobs.view(-1, lprobs.size(-1)), target.view(-1))

                # logging
                accum_loss += loss.item()
                # back-prop
                scaler.scale(loss).backward()

            scaler.unscale_(optimizer)
            Loading [MathJax]/jax/output/HTML-CSS/fonts/TeX/fontdata.js] (sample_size or 1.0)) # (sample_size or 1.0) handles the case of a zero gradient
            gnorm = nn.utils.clip_grad_norm_(model.parameters(), config.clip_norm) # 梯度裁剪 防止梯度爆炸
```

```

scaler.step(optimizer)
scaler.update()

# logging
loss_print = accum_loss/sample_size
stats["loss"].append(loss_print)
progress.set_postfix(loss=loss_print)
if config.use_wandb:
    wandb.log({
        "train/loss": loss_print,
        "train/grad_norm": gnorm.item(),
        "train/lr": optimizer.rate(),
        "train/sample_size": sample_size,
    })

loss_print = np.mean(stats["loss"])
logger.info(f"training loss: {loss_print:.4f}")
return stats

```

## Validation & Inference 檢驗和推論

為防止訓練發生過度擬合，每過一段時間要做一次檢測，計算模型在未看過的資料上的表現。

- 過程基本上和training一樣，另外加上 inference
- 檢驗完畢可順便儲存模型參數

單看 validation loss，我們很難知道模型真實的效能

- 直接用當前模型去生成翻譯結果 (hypothesis)，再和正確答案 (reference) 計算 BLEU score
- 也可用肉眼看翻譯結果的好壞
- 我們用 fairseq 寫好的 sequence generator 來進行 beam search 生成翻譯結果

Loading [MathJax]/jax/output/HTML-CSS/fonts/TeX/fontdata.js

```
In [ ]: # fairseq 的 beam search generator
# 給定模型和輸入序列，用 beam search 生成翻譯結果
sequence_generator = task.build_generator([model], config)

def decode(toks, dictionary):
    # 從 Tensor 轉成人看得懂的句子
    s = dictionary.string(
        toks.int().cpu(),
        config.post_process,
    )
    return s if s else "<unk>"

def inference_step(sample, model):
    gen_out = sequence_generator.generate([model], sample)
    srcs = []
    hyps = []
    refs = []
    for i in range(len(gen_out)):
        # 對於每個 sample，收集輸入，輸出和參考答案，稍後計算 BLEU
        srcs.append(decode(
            utils.strip_pad(sample["net_input"]["src_tokens"][i], task.source_dictionary.pad()),
            task.source_dictionary,
        ))
        hyps.append(decode(
            gen_out[i][0]["tokens"], # 0 代表取出 beam 內分數第一的輸出結果
            task.target_dictionary,
        ))
        refs.append(decode(
            utils.strip_pad(sample["target"][i], task.target_dictionary.pad()),
            task.target_dictionary,
        ))
    return srcs, hyps, refs
```

Loading [MathJax]/jax/output/HTML-CSS/fonts/TeX/fontdata.js

```
In [ ]: import shutil
import sacrebleu

def validate(model, task, criterion, log_to_wandb=True):
    logger.info('begin validation')
    itr = load_data_iterator(task, "valid", 1, config.max_tokens, config.num_workers).next_epoch_itr(shuffle=False)

    stats = {"loss":[], "bleu": 0, "srcs":[], "hyps":[], "refs":[]}
    srcs = []
    hyps = []
    refs = []

    model.eval()
    progress = tqdm.tqdm(itr, desc=f"validation", leave=False)
    with torch.no_grad():
        for i, sample in enumerate(progress):
            # validation loss
            sample = utils.move_to_cuda(sample, device=device)
            net_output = model.forward(**sample["net_input"])

            lprobs = F.log_softmax(net_output[0], -1)
            target = sample["target"]
            sample_size = sample["ntokens"]
            loss = criterion(lprobs.view(-1, lprobs.size(-1)), target.view(-1)) / sample_size
            progress.set_postfix(valid_loss=loss.item())
            stats["loss"].append(loss)

            # 進行推論
            s, h, r = inference_step(sample, model)
            srcs.extend(s)
            hyps.extend(h)
            refs.extend(r)

            tok = 'zh' if task.cfg.target_lang == 'zh' else '13a'
            stats["loss"] = torch.stack(stats["loss"]).mean().item()
            stats["bleu"] = sacrebleu.corpus_bleu(hyps, [refs], tokenize=tok) # 計算BLEU score
            stats["srcs"] = srcs
            stats["hyps"] = hyps
            stats["refs"] = refs
```

Loading [MathJax]/jax/output/HTML-CSS/fonts/TeX/fontdata.js  
if config.use\_wandb and log\_to\_wandb:

```
wandb.log({  
    "valid/loss": stats["loss"],  
    "valid/bleu": stats["bleu"].score,  
}, commit=False)  
  
showid = np.random.randint(len(hyps))  
logger.info("example source: " + srcs[showid])  
logger.info("example hypothesis: " + hyps[showid])  
logger.info("example reference: " + refs[showid])  
  
# show bleu results  
logger.info(f"validation loss:\t{stats['loss']:.4f}")  
logger.info(stats["bleu"].format())  
return stats
```

## 儲存及載入模型參數

Loading [MathJax]/jax/output/HTML-CSS/fonts/TeX/fontdata.js

```
In [ ]: def validate_and_save(model, task, criterion, optimizer, epoch, save=True):
    stats = validate(model, task, criterion)
    bleu = stats['bleu']
    loss = stats['loss']
    if save:
        # save epoch checkpoints
        savedir = Path(config.savedir).absolute()
        savedir.mkdir(parents=True, exist_ok=True)

        check = {
            "model": model.state_dict(),
            "stats": {"bleu": bleu.score, "loss": loss},
            "optim": {"step": optimizer._step}
        }
        torch.save(check, savedir/f"checkpoint{epoch}.pt")
        shutil.copy(savedir/f"checkpoint{epoch}.pt", savedir/f"checkpoint_last.pt")
        logger.info(f"saved epoch checkpoint: {savedir}/checkpoint{epoch}.pt")

        # save epoch samples
        with open(savedir/f"samples{epoch}. {config.source_lang}-{config.target_lang}.txt", "w") as f:
            for s, h in zip(stats["srcs"], stats["hypss"]):
                f.write(f"{s}\t{h}\n")

        # get best valid bleu
        if getattr(validate_and_save, "best_bleu", 0) < bleu.score:
            validate_and_save.best_bleu = bleu.score
            torch.save(check, savedir/f"checkpoint_best.pt")

        del_file = savedir / f"checkpoint{epoch - config.keep_last_epochs}.pt"
        if del_file.exists():
            del_file.unlink()

    return stats

def try_load_checkpoint(model, optimizer=None, name=None):
    name = name if name else "checkpoint_last.pt"
    checkpath = Path(config.savedir)/name
    if checkpath.exists():
        check = torch.load(checkpath)
        model.load_state_dict(check["model"])
        stats = check["stats"]
```

Loading [MathJax]/jax/output/html-CSS/fonts/TeX/fontdata.js  
 stats = check["stats"]

```

step = "unknown"
if optimizer != None:
    optimizer._step = step = check["optim"]["step"]
logger.info(f"loaded checkpoint {checkpath}: step={step} loss={stats['loss']} bleu={stats['bleu']} ")
else:
    logger.info(f"no checkpoints found at {checkpath}!")

```

# 主程式

## 訓練迴圈

In [ ]: model = model.to(device=device)  
criterion = criterion.to(device=device)

In [ ]: !nvidia-smi

In [ ]: logger.info("task: {}".format(task.\_\_class\_\_.\_\_name\_\_))
logger.info("encoder: {}".format(model.encoder.\_\_class\_\_.\_\_name\_\_))
logger.info("decoder: {}".format(model.decoder.\_\_class\_\_.\_\_name\_\_))
logger.info("criterion: {}".format(criterion.\_\_class\_\_.\_\_name\_\_))
logger.info("optimizer: {}".format(optimizer.\_\_class\_\_.\_\_name\_\_))
logger.info(
 "num. model params: {}, (num. trained: {},)".format(
 sum(p.numel() for p in model.parameters()),
 sum(p.numel() for p in model.parameters() if p.requires\_grad),
 )
)
logger.info(f"max tokens per batch = {config.max\_tokens}, accumulate steps = {config.accum\_steps}")

Loading [MathJax]/jax/output/HTML-CSS/fonts/TeX/fontdata.js

```
In [ ]: epoch_itr = load_data_iterator(task, "train", config.start_epoch, config.max_tokens, config.num_workers)
try_load_checkpoint(model, optimizer, name=config.resume)
while epoch_itr.next_epoch_idx <= config.max_epoch:
    # train for one epoch
    train_one_epoch(epoch_itr, model, task, criterion, optimizer, config.accum_steps)
    stats = validate_and_save(model, task, criterion, optimizer, epoch=epoch_itr.epoch)
    logger.info("end of epoch {}".format(epoch_itr.epoch))
    epoch_itr = load_data_iterator(task, "train", epoch_itr.next_epoch_idx, config.max_tokens, config.num_workers)
```

## Submission 繳交檔案

```
In [ ]: # 把幾個 checkpoint 平均起來可以達到 ensemble 的效果
checkdir=config.savedir
!python ./fairseq/scripts/average_checkpoints.py \
--inputs {checkdir} \
--num-epoch-checkpoints 5 \
--output {checkdir}/avg_last_5_checkpoint.pt
```

## 確認生成繳交檔案的模型參數

```
In [ ]: # checkpoint_last.pt : 最後一次檢驗的檔案
# checkpoint_best.pt : 檢驗 BLEU 最高的檔案
# avg_last_5_checkpoint.pt: 最5後個檔案平均
try_load_checkpoint(model, name="checkpoint_best.pt")
validate(model, task, criterion, log_to_wandb=False)
None
```

## 進行預測

Loading [MathJax]/jax/output/HTML-CSS/fonts/TeX/fontdata.js

```
In [ ]: def generate_prediction(model, task, split="test", outfile="../prediction.txt"):
    task.load_dataset(split=split, epoch=1)
    itr = load_data_iterator(task, split, 1, config.max_tokens, config.num_workers).next_epoch_itr(shuffle=False)

    idxs = []
    hyps = []

    model.eval()
    progress = tqdm.tqdm(itr, desc=f"prediction")
    with torch.no_grad():
        for i, sample in enumerate(progress):
            # validation loss
            sample = utils.move_to_cuda(sample, device=device)

            # 進行推論
            s, h, r = inference_step(sample, model)

            hyps.append(h)
            idxs.append(list(sample['id']))

    # 根據 preprocess 時的順序排列
    hyps = [x for _, x in sorted(zip(idxs, hyps))]

    with open(outfile, "w") as f:
        for h in hyps:
            f.write(h + "\n")
```

```
In [ ]: generate_prediction(model, task)
```

```
In [ ]: raise
```

## Back-translation

訓練一個反向的翻譯模型

1. 將實驗的參數設定表中(config)的source\_lang與target\_lang互相交換
2. 將實驗的參數設定表中(config)的savedir更改(ex. "./checkpoints/rnn-back")
3. 訓練一個反向模型

## 利用反向模型生成額外資料

### 下載 monolingual data

```
In [ ]: mono_dataset_name = 'mono'
```

Loading [MathJax]/jax/output/HTML-CSS/fonts/TeX/fontdata.js

```
In [ ]: mono_prefix = Path(data_dir).absolute() / mono_dataset_name
mono_prefix.mkdir(parents=True, exist_ok=True)

urls = (
    "https://onedrive.live.com/download?cid=3E549F3B24B238B4&resid=3E549F3B24B238B4%214986&authkey=AANUKbGfZx0kM80",
    # # If the above links die, use the following instead.
    "# https://www.csie.ntu.edu.tw/~r09922057/ML2021-hw5/ted_zh_corpus.deduped.gz",
    # # If the above links die, use the following instead.
    "# https://mega.nz/#!vMNnDShR!4eHDxzlpzIpdpeQTD-htatU_C7QwcBTwGDaSeBqH534",
)
file_names = (
    'ted_zh_corpus.deduped.gz',
)

for u, f in zip(urls, file_names):
    path = mono_prefix/f
    if not path.exists():
        if 'mega' in u:
            !megadl {u} --path {path}
        else:
            !wget {u} -O {path}
    else:
        print(f'{f} is exist, skip downloading')
    if path.suffix == ".tgz":
        !tar -xvf {path} -C {prefix}
    elif path.suffix == ".zip":
        !unzip -o {path} -d {prefix}
    elif path.suffix == ".gz":
        !gzip -fk {path}
```

## TODO: 清理資料集

1. 將太長、太短的句子移除
2. 統一標點符號

hint: 可以使用clean\_s()來協助

Loading [MathJax]/jax/output/HTML-CSS/fonts/TeX/fontdata.js

In [ ]:

## TODO: Subword Units

用反向模型的 spm model 將資料切成 subword units

hint: spm model 的路徑為 DATA/raw-data/[dataset]/spm[vocab\_num].model

In [ ]:

## Binarize

使用fairseq將資料轉為binary

```
In [ ]: binpath = Path('./DATA/data-bin', mono_dataset_name)
src_dict_file = './DATA/data-bin/ted2020/dict.en.txt'
tgt_dict_file = src_dict_file
monopref = str(mono_prefix/"mono.tok") # whatever filepath you get after applying subword tokenization
if binpath.exists():
    print(binpath, "exists, will not overwrite!")
else:
    !python -m fairseq_cli.preprocess
        --source-lang 'zh' \
        --target-lang 'en' \
        --trainpref {monopref} \
        --destdir {binpath} \
        --srctext {src_dict_file} \
        --tgtdict {tgt_dict_file} \
        --workers 2
```

Loading [MathJax]/jax/output/HTML-CSS/fonts/TeX/fontdata.js

## TODO: 生成反向翻譯資料

將 binarized data 加入原本的資料夾中並用一個 split\_name 取名

ex. ./DATA/data-bin/ted2020/[split\_name].zh-en.[{"en", "zh"}].[{"bin", "idx"}]

便可以使用 generate\_prediction(model, task, split="split\_name")來產生翻譯資料

```
In [ ]: # 將 binarized data 加入原本的資料夾中並用一個 split_name 取名
# ex. ./DATA/data-bin/ted2020/\[split_name\].zh-en.\["en", "zh"\].\[{"bin", "idx"}\]
!cp ./DATA/data-bin/mono/train.zh-en.zh.bin ./DATA/data-bin/ted2020/mono.zh-en.zh.bin
!cp ./DATA/data-bin/mono/train.zh-en.zh.idx ./DATA/data-bin/ted2020/mono.zh-en.zh.idx
!cp ./DATA/data-bin/mono/train.zh-en.en.bin ./DATA/data-bin/ted2020/mono.zh-en.en.bin
!cp ./DATA/data-bin/mono/train.zh-en.en.idx ./DATA/data-bin/ted2020/mono.zh-en.en.idx
```

```
In [ ]: # hint: 用反向模型在 split='mono' 上進行預測, 生成 prediction_file
# generate_prediction( ... ,split=... ,outfile=... )
```

## TODO: 產生新的dataset

1. 將翻譯出來的資料與原先的訓練資料結合
2. 使用之前的spm model切出成Subword Units
3. 重新使用fairseq將資料轉為binary

Loading [MathJax]/jax/output/HTML-CSS/fonts/TeX/fontdata.js

```
In [ ]: # 合併剛剛生成的 prediction_file (.en) 以及中文 mono.zh (.zh)
#
# hint: 在此用剛剛的 spm model 對 prediction_file 進行切斷詞
# spm_model.encode(line, out_type=str)
# output: ./DATA/rawdata/mono/mono.tok.en & mono.tok.zh
#
# hint: 在此用 fairseq 把這些檔案再 binarize
# binpath = Path('./DATA/data-bin/synthetic')
# src_dict_file = './DATA/data-bin/ted2020/dict.en.txt'
# tgt_dict_file = src_dict_file
# monopref = './DATA/rawdata/mono/mono.tok' # or whatever path after applying subword tokenization, w/o the suffix (.zh/.en)
# if binpath.exists():
#     print(binpath, "exists, will not overwrite!")
# else:
#     !python -m fairseq_cli.preprocess \
#         --source-lang 'zh' \
#         --target-lang 'en' \
#         --trainpref {monopref} \
#         --destdir {binpath} \
#         --srcdict {src_dict_file} \
#         --tgtdict {tgt_dict_file} \
#         --workers 2
```

```
In [ ]: # 這裡用剛剛準備的檔案合併原先 ted2020 來生成最終 back-translation 的資料
!cp -r ./DATA/data-bin/ted2020/ ./DATA/data-bin/ted2020_with_mono/
!cp ./DATA/data-bin/synthetic/train.zh-en.zh.bin ./DATA/data-bin/ted2020_with_mono/train1.en-zh.zh.bin
!cp ./DATA/data-bin/synthetic/train.zh-en.zh.idx ./DATA/data-bin/ted2020_with_mono/train1.en-zh.zh.idx
!cp ./DATA/data-bin/synthetic/train.zh-en.en.bin ./DATA/data-bin/ted2020_with_mono/train1.en-zh.en.bin
!cp ./DATA/data-bin/synthetic/train.zh-en.en.idx ./DATA/data-bin/ted2020_with_mono/train1.en-zh.en.idx
```

## TODO: 重新訓練

當已經產生新的資料集

1. 將實驗的參數設定表(config)中的datadir改為新的資料集("./DATA/data-bin/ted2020\_with\_mono")
2. 將實驗的參數設定表(config)中的source\_lang與target\_lang設定還原("en", "zh")
3. 將實驗的參數設定表(config)中的savedir更改(ex. "./checkpoints/rnn-bt")
4. 重新訓練

## References

1. Ott, M., Edunov, S., Baevski, A., Fan, A., Gross, S., Ng, N., ... & Auli, M. (2019, June). fairseq: A Fast, Extensible Toolkit for Sequence Modeling. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations) (pp. 48-53).
2. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017, December). Attention is all you need. In Proceedings of the 31st International Conference on Neural Information Processing Systems (pp. 6000-6010).
3. Reimers, N., & Gurevych, I. (2020, November). Making Monolingual Sentence Embeddings Multilingual Using Knowledge Distillation. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP) (pp. 4512-4525).
4. Tiedemann, J. (2012, May). Parallel Data, Tools and Interfaces in OPUS. In Lrec (Vol. 2012, pp. 2214-2218).
5. Kudo, T., & Richardson, J. (2018, November). SentencePiece: A simple and language independent subword tokenizer and detokenizer for Neural Text Processing. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations (pp. 66-71).
6. Sennrich, R., Haddow, B., & Birch, A. (2016, August). Improving Neural Machine Translation Models with Monolingual Data. In Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers) (pp. 86-96).
7. Edunov, S., Ott, M., Auli, M., & Grangier, D. (2018). Understanding Back-Translation at Scale. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing (pp. 489-500).
8. <https://github.com/ajinkyakulkarni14/TED-Multilingual-Parallel-Corpus> (<https://github.com/ajinkyakulkarni14/TED-Multilingual-Parallel-Corpus>)
9. <https://ithelp.ithome.com.tw/articles/10233122> (<https://ithelp.ithome.com.tw/articles/10233122>)
10. <https://nlp.seas.harvard.edu/2018/04/03/attention.html> (<https://nlp.seas.harvard.edu/2018/04/03/attention.html>)

In [ ]: Loading [MathJax]/jax/output/HTML-CSS/fonts/TeX/fontdata.js