TRANSLATION WITH A SEQUENCE TO SEQUENCE NETWORK AND ATTENTION

- Translate from French to English
- Sequence to sequence network
  - : two RNN work together to transform one sequence to another.
    - · Encoder: condences an input sequence into a vector
    - · Decoder: unfolds the vector into a new sequence

## Requirements

from \_\_future \_\_ import unicode\_literals, print \_ function, division

from io import open

import unicodedata

import string

import re

import random

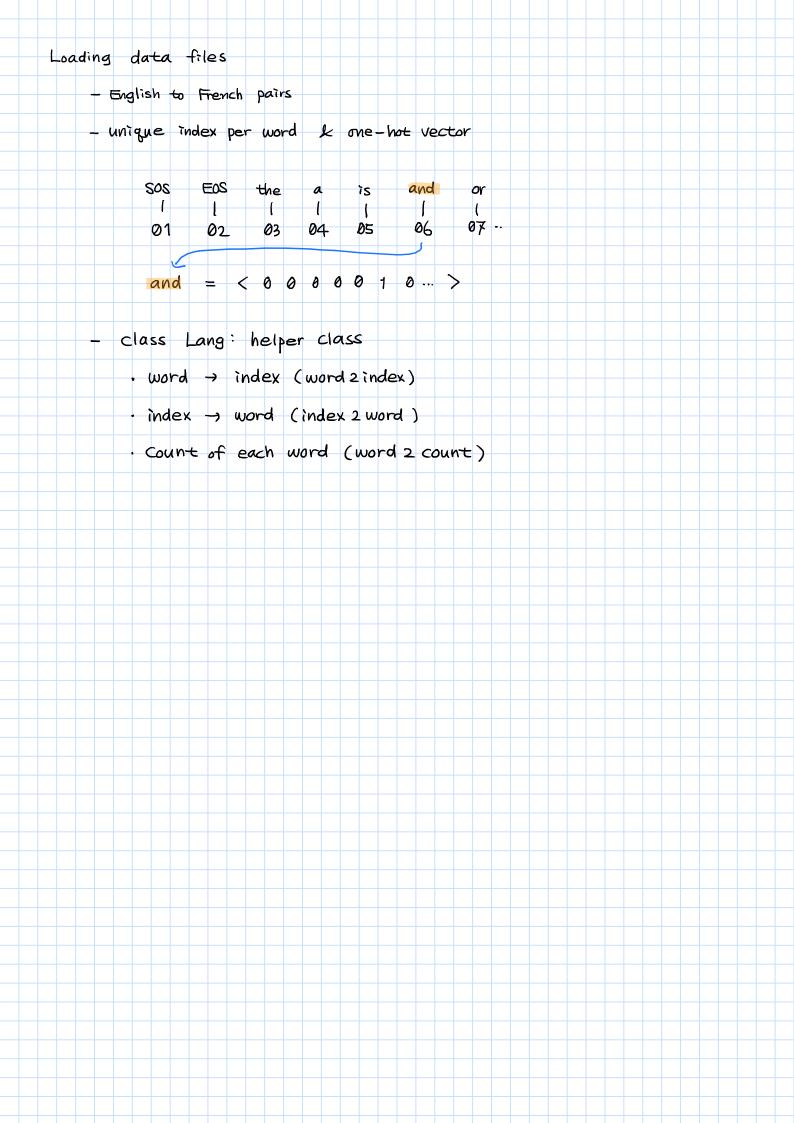
import torch

from torch nn import nn

from torch import optim

import torch nn. functional as F

device = torch device ("cuda" if torch cuda is available () else "cpu")



```
SOS_token = 0
EOS_token = 1
class Lang:
    def __init __ (self, name):
        self. name = name
        self. word 2 index = {}
        Self. word 2 count = { }
        self. index 2 word = { 0: "SOS", 1: "EOS" }
        self. n_words = 2 # Count SOS and EOS
    def add Sentence (self, sentence):
        for word in sentence. split (' ');
            self. add Word (word)
    def add Word (self, word):
        if word not in self. word 2 index:
            self. word 2 index [word] = self. n_words
            self. word 2 Count [word] = 1
            self. index 2 word [self. n_words] = word
            self. n-words += 1
        else :
            self. word 2 count [word] += 1
```

```
To simplify:
- turn Unicode characters to ASCII
- make everything lowercase
- trim most punctuation
def unicode To Ascii (s):
    return ". join (
         c for c in unicodedata normalize ('NFD', s)
        if unicodedata category (c) != 'Mn'
# Lowercase, trim, and remove non-letter characters
def normalize String (s):
    S = unicode To Ascii (s. lower (). strip())
    s = re.sub(r"([.!?])", r" \setminus 1", s)
    S = re.sub(r"[^a-zA-z.!?]+", r" ", s)
    return s
```

```
To read the data file:
- split the file into lines
- split lines into pairs
To translate from other language -> English: reverse flag
def read Largs (larg1, larg1, reverse = False):
     print ("Reading lines ...")
    # Read the file and split into lines
    lines = open ('data/%s-%s.txt' % (lang1, lang2), encoding = 'utf-8').
         read() strip(). split('\n')
    # split every line into pairs and normalize
     pairs = [[normalizeString(s) for s in 1.split("It")] for 1 in lines]
     # Reverse pairs, make Lang instances
     if reverse:
         pairs = [list (reversed (p)) for p in pairs]
        input_lang = Lang (lang 2)
         output _ lang = Lang (lang 1)
    else :
         input _ lang = Lang (lang 1)
         Output _ lang = Lang (lang 2)
    return input - lang, output - lang, pairs
```

```
To train quickly:
- trim the dataset to only relatively short & simple sentences.
- maximum length: 10 words (includes ending punctuation)
- filtering to sentences that translate to the form "I am" or "He is" etc.
MAX_ LENGTH = 10
eng_prefixes = (
   "i am", "i am", "he is", "he s",
   "she is", "she s", "you are", "you re",
   "we are", "we re", "they are", "they re'
def firterPair (p):
    return (en(p[0].split('')) < MAX_LENGTH and \
        len(p[1] split(' ')) < MAX_LENGH and \
        p[1] startswith (eng_prefixes)
def fitter Pairs (pairs):
             [pair for pair in pairs if fitterPair (pair)]
```

Full process for preparing the data is: - Read text file and split into lines, split lines into pairs - Normalize text, filter by length & content - Make word lists from sentences in pairs