

Lab Week 9: Chatbots 1 – Data preparation



Conversational models are a hot topic in artificial intelligence research. Chatbots can be found in a variety of settings, including customer service applications and online helpdesks. These bots are often powered by retrieval-based models, which output predefined responses to questions of certain forms. In this Lab 9 and 10 we will train a simple chatbot using movie scripts from the [Cornell Movie-Dialogs Corpus](#).

We focus in this Lab 9 on Data preparation.

This lab is based on the Pytorch [ChatBot Tutorial](#) and borrows code from the following sources:

1. Yuan-Kuei Wu's pytorch-chatbot implementation: <https://github.com/ywk991112/pytorch-chatbot>
2. Sean Robertson's practical-pytorch seq2seq-translation example: <https://github.com/spro/practical-pytorch/tree/master/seq2seq-translation>
3. FloydHub's Cornell Movie Corpus preprocessing code: <https://github.com/floydhub/textutil-preprocess-cornell-movie-corpus>

To start, you need to download the data ZIP file from Canvas or using the below command and put in a `data/` directory under your current working directory.

```
wget -c http://www.mpi-sws.org/~cristian/data/cornell\_movie\_dialogs\_corpus.zip
```

The [Cornell Movie-Dialogs Corpus](#) is a rich dataset of movie character dialog:

- 220,579 conversational exchanges between 10,292 pairs of movie characters
- 9,035 characters from 617 movies
- 304,713 total utterances

This dataset is large and diverse, and there is a great variation of language formality, time periods, sentiment, etc. Our hope is that this diversity makes our model robust to many forms of inputs and queries.

Now, let's import some necessities.

```
from __future__ import absolute_import
from __future__ import division
from __future__ import print_function
from __future__ import unicode_literals
```

```
import torch
from torch.jit import script, trace
import torch.nn as nn
from torch import optim
import torch.nn.functional as F
import csv
import random
import re
import os
import unicodedata
import codecs
from io import open
import itertools
import math
```

```
USE_CUDA = torch.cuda.is_available()
device = torch.device("cuda" if USE_CUDA else "cpu")
```

Load & Preprocess Data

You need to take a look at some lines of your datafile to see the original format.

```
corpus_name = "cornell movie-dialogs corpus"
corpus = os.path.join("data", corpus_name)

def printLines(file, n=10):
    with open(file, 'rb') as datafile:
        lines = datafile.readlines()
        for line in lines[:n]:
            print(line)

printLines(os.path.join(corpus, "movie_lines.txt"))
```

Create formatted data file

For convenience, we'll create a nicely formatted data file in which each line contains a tab-separated *query sentence* and a *response sentence* pair.

The following functions facilitate the parsing of the raw *movie_lines.txt* data file.

- `loadLines` splits each line of the file into a dictionary of fields (lineID, characterID, movieID, character, text)

- `loadConversations` groups fields of lines from `loadLines` into conversations based on `movie_conversations.txt`
- `extractSentencePairs` extracts pairs of sentences from conversations

Splits each line of the file into a dictionary of fields

`def loadLines(fileName, fields):`

```
    lines = {}
    with open(fileName, 'r', encoding='iso-8859-1') as f:
        for line in f:
            values = line.split(" +++$+++ ")
            # Extract fields
            lineObj = {}
            for i, field in enumerate(fields):
                lineObj[field] = values[i]
            lines[lineObj['lineID']] = lineObj
    return lines
```

*# Groups fields of lines from `loadLines` into conversations based on *movie_conversations.txt**

`def loadConversations(fileName, lines, fields):`

```
    conversations = []
    with open(fileName, 'r', encoding='iso-8859-1') as f:
        for line in f:
            values = line.split(" +++$+++ ")
            # Extract fields
            convObj = {}
            for i, field in enumerate(fields):
                convObj[field] = values[i]
            # Convert string to list (convObj["utteranceIDs"] == "['L598485', 'L598486', ...]")
            utterance_id_pattern = re.compile('L[0-9]+')
            lineIds = utterance_id_pattern.findall(convObj["utteranceIDs"])
            # Reassemble lines
            convObj["lines"] = []
            for lineId in lineIds:
                convObj["lines"].append(lines[lineId])
            conversations.append(convObj)
    return conversations
```

Extracts pairs of sentences from conversations

`def extractSentencePairs(conversations):`

```
    qa_pairs = []
    for conversation in conversations:
        # Iterate over all the lines of the conversation
        for i in range(len(conversation["lines"]) - 1): # We ignore the last line (no answer for it)
            inputLine = conversation["lines"][i]["text"].strip()
            targetLine = conversation["lines"][i+1]["text"].strip()
            # Filter wrong samples (if one of the lists is empty)
            if inputLine and targetLine:
                qa_pairs.append([inputLine, targetLine])
    return qa_pairs
```

Now we'll call these functions and create the file. We'll call it *formatted_movie_lines.txt*.

```

# Define path to new file
datafile = os.path.join(corpus, "formatted_movie_lines.txt")

delimiter = '\t'
# Unescape the delimiter
delimiter = str(codecs.decode(delimiter, "unicode_escape"))

# Initialize lines dict, conversations list, and field ids
lines = {}
conversations = []
MOVIE_LINES_FIELDS = ["lineID", "characterID", "movieID", "character", "text"]
MOVIE_CONVERSATIONS_FIELDS = ["character1ID", "character2ID", "movieID", "utteranceIDs"]

# Load lines and process conversations
print("\nProcessing corpus...")
lines = loadLines(os.path.join(corpus, "movie_lines.txt"), MOVIE_LINES_FIELDS)
print("\nLoading conversations...")
conversations = loadConversations(os.path.join(corpus, "movie_conversations.txt"),
                                   lines, MOVIE_CONVERSATIONS_FIELDS)

# Write new csv file
print("\nWriting newly formatted file...")
with open(datafile, 'w', encoding='utf-8') as outputfile:
    writer = csv.writer(outputfile, delimiter=delimiter, lineterminator='\n')
    for pair in extractSentencePairs(conversations):
        writer.writerow(pair)

# Print a sample of lines
print("\nSample lines from file:")
printLines(datafile)

```

Load and trim data

Our next step is to create a **vocabulary** and load query/response sentence pairs into memory.

Note that we are dealing with sequences of **words**, which do not have an implicit mapping to a discrete numerical space. Thus, we must create one by mapping each unique word that we encounter in our dataset to an index value.

For this we define a `Voc` class, which keeps a mapping from words to indexes, a reverse mapping of indexes to words, a count of each word and a total word count. The class provides methods for adding a word to the vocabulary (`addWord`), adding all words in a sentence (`addSentence`) and trimming infrequently seen words (`trim`). More on trimming later.

```

# Default word tokens
PAD_token = 0 # Used for padding short sentences
SOS_token = 1 # Start-of-sentence token
EOS_token = 2 # End-of-sentence token

class Voc:
    def __init__(self, name):
        self.name = name
        self.trimmed = False

```

```

self.word2index = {}
self.word2count = {}
self.index2word = {PAD_token: "PAD", SOS_token: "SOS", EOS_token: "EOS"}
self.num_words = 3 # Count SOS, EOS, PAD

def addSentence(self, sentence):
    for word in sentence.split(' '):
        self.addWord(word)

def addWord(self, word):
    if word not in self.word2index:
        self.word2index[word] = self.num_words
        self.word2count[word] = 1
        self.index2word[self.num_words] = word
        self.num_words += 1
    else:
        self.word2count[word] += 1

# Remove words below a certain count threshold
def trim(self, min_count):
    if self.trimmed:
        return
    self.trimmed = True

    keep_words = []

    for k, v in self.word2count.items():
        if v >= min_count:
            keep_words.append(k)

    print('keep_words {} / {} = {:.4f}'.format(
        len(keep_words), len(self.word2index), len(keep_words) / len(self.word2index)
    ))

# Reinitialize dictionaries
self.word2index = {}
self.word2count = {}
self.index2word = {PAD_token: "PAD", SOS_token: "SOS", EOS_token: "EOS"}
self.num_words = 3 # Count default tokens

for word in keep_words:
    self.addWord(word)

```

Now we can assemble our vocabulary and query/response sentence pairs. Before we are ready to use this data, we must perform some preprocessing.

First, we must convert the Unicode strings to ASCII using `unicodeToAscii`. Next, we should convert all letters to lowercase and trim all non-letter characters except for basic punctuation (`normalizeString`). Finally, to aid in training convergence, we will filter out sentences with length greater than the `MAX_LENGTH` threshold (`filterPairs`).

`MAX_LENGTH = 10` # Maximum sentence length to consider

```

# Turn a Unicode string to plain ASCII, thanks to
# https://stackoverflow.com/a/518232/2809427
def unicodeToAscii(s):
    return ''.join(

```

```

        c for c in unicodedata.normalize('NFD', s)
        if unicodedata.category(c) != 'Mn'
    )

# Lowercase, trim, and remove non-letter characters
def normalizeString(s):
    s = unicodeToAscii(s.lower().strip())
    s = re.sub(r"([.!?])", r" \1", s)
    s = re.sub(r"^a-zA-Z.!?]+", r" ", s)
    s = re.sub(r"\s+", r" ", s).strip()
    return s

# Read query/response pairs and return a voc object
def readVocs(datafile, corpus_name):
    print("Reading lines...")
    # Read the file and split into lines
    lines = open(datafile, encoding='utf-8').\
        read().strip().split('\n')
    # Split every line into pairs and normalize
    pairs = [[normalizeString(s) for s in l.split('\t')] for l in lines]
    voc = Voc(corpus_name)
    return voc, pairs

# Returns True iff both sentences in a pair 'p' are under the MAX_LENGTH threshold
def filterPair(p):
    # Input sequences need to preserve the last word for EOS token
    return len(p[0].split(' ')) < MAX_LENGTH and len(p[1].split(' ')) < MAX_LENGTH

# Filter pairs using filterPair condition
def filterPairs(pairs):
    return [pair for pair in pairs if filterPair(pair)]

# Using the functions defined above, return a populated voc object and pairs list
def loadPrepareData(corpus, corpus_name, datafile, save_dir):
    print("Start preparing training data ...")
    voc, pairs = readVocs(datafile, corpus_name)
    print("Read {!s} sentence pairs".format(len(pairs)))
    pairs = filterPairs(pairs)
    print("Trimmed to {!s} sentence pairs".format(len(pairs)))
    print("Counting words...")
    for pair in pairs:
        voc.addSentence(pair[0])
        voc.addSentence(pair[1])
    print("Counted words:", voc.num_words)
    return voc, pairs

# Load/Assemble voc and pairs
save_dir = os.path.join("data", "save")
voc, pairs = loadPrepareData(corpus, corpus_name, datafile, save_dir)
# Print some pairs to validate
print("\npairs:")
for pair in pairs[:10]:
    print(pair)

```

Another technique that is beneficial to achieving faster convergence during training is trimming rarely used words out of our vocabulary. Decreasing the feature space will also soften the difficulty of the function that the model must learn to approximate. We will do this as a two-step process:

1. Trim words used under `MIN_COUNT` threshold using the `voc.trim` function.
2. Filter out pairs with trimmed words.

`MIN_COUNT = 3` # Minimum word count threshold for trimming

```
def trimRareWords(voc, pairs, MIN_COUNT):
    # Trim words used under the MIN_COUNT from the voc
    voc.trim(MIN_COUNT)
    # Filter out pairs with trimmed words
    keep_pairs = []
    for pair in pairs:
        input_sentence = pair[0]
        output_sentence = pair[1]
        keep_input = True
        keep_output = True
        # Check input sentence
        for word in input_sentence.split(' '):
            if word not in voc.word2index:
                keep_input = False
                break
        # Check output sentence
        for word in output_sentence.split(' '):
            if word not in voc.word2index:
                keep_output = False
                break

        # Only keep pairs that do not contain trimmed word(s) in their input or output sentence
        if keep_input and keep_output:
            keep_pairs.append(pair)

    print('Trimmed from {} pairs to {}, {:.4f} of total'.format(len(pairs), len(keep_pairs), len(keep_pairs) /
len(pairs)))
    return keep_pairs

# Trim voc and pairs
pairs = trimRareWords(voc, pairs, MIN_COUNT)
```

Although we have put a great deal of effort into preparing and massaging our data into a nice vocabulary object and list of sentence pairs, our models will ultimately expect numerical torch tensors as inputs. One way to prepare the processed data for the models is to use a batch size of 1, meaning that all we have to do is convert the words in our sentence pairs to their corresponding indexes from the vocabulary and feed this later (In Lab 10) to the models.