

TOPIC OF STUDY: MACHINE TRANSLATION

Code:

```
!pip install transformers sentencepiece datasets

from datasets import load_dataset
from google.colab import drive
from IPython.display import display
from IPython.html import widgets
import matplotlib.pyplot as plt
import numpy as np
import seaborn as sns
import torch
from torch import optim
from torch.nn import functional as F
from transformers import AdamW, AutoModelForSeq2SeqLM, AutoTokenizer
from transformers import get_linear_schedule_with_warmup
from tqdm import tqdm_notebook

sns.set()

drive.mount('/content/gdrive')

# Use 'google/mt5-small' for non-pro cloab users
model_repo = 'google/mt5-base'
model_path = '/content/gdrive/My Drive/mt5_translation.pt'
max_seq_len = 20

"""# Load Tokenizer & Model"""

tokenizer = AutoTokenizer.from_pretrained(model_repo)

# Model description: https://huggingface.co/google/mt5-base
model = AutoModelForSeq2SeqLM.from_pretrained(model_repo)
model = model.cuda()

"""# Overview and Quick Test"""

token_ids = tokenizer.encode(
    '<jp> This will be translated to Japanese! (hopefully)',
    return_tensors='pt').cuda()
print(token_ids)

model_out = model.generate(token_ids)
print(model_out)
```

```

output_text = tokenizer.convert_tokens_to_string(
    tokenizer.convert_ids_to_tokens(model_out[0]))
print(output_text)

"""# Steps
1. Load the pretrained model and tokenizer
2. Load dataset
3. Transform dataset into input (entails a minor model change)
4. Train/finetune the model on our dataset
5. Test the model

# Test Tokenizer
"""

example_input_str = '<jp> This is just a test nbuig.'
# example_input_str = 'これは普通のテスト'
input_ids = tokenizer.encode(example_input_str, return_tensors='pt')
print('Input IDs:', input_ids)

tokens = tokenizer.convert_ids_to_tokens(input_ids[0])
print('Tokens:', tokens)

sorted(tokenizer.vocab.items(), key=lambda x: x[1])

"""# Prepare Dataset"""

# Source: https://huggingface.co/datasets/alt
dataset = load_dataset('alt')

train_dataset = dataset['train']
test_dataset = dataset['test']

train_dataset[0]

LANG_TOKEN_MAPPING = {
    'en': '<en>',
    'ja': '<jp>',
    'zh': '<zh>'
}

special_tokens_dict = {'additional_special_tokens':
list(LANG_TOKEN_MAPPING.values())}
tokenizer.add_special_tokens(special_tokens_dict)
model.resize_token_embeddings(len(tokenizer))

token_ids = tokenizer.encode(
    example_input_str, return_tensors='pt', padding='max_length',
    truncation=True, max_length=max_seq_len)

```

```

print(token_ids)

tokens = tokenizer.convert_ids_to_tokens(token_ids[0])
print(tokens)

def encode_input_str(text, target_lang, tokenizer, seq_len,
                    lang_token_map=LANG_TOKEN_MAPPING):
    target_lang_token = lang_token_map[target_lang]

    # Tokenize and add special tokens
    input_ids = tokenizer.encode(
        text = target_lang_token + text,
        return_tensors = 'pt',
        padding = 'max_length',
        truncation = True,
        max_length = seq_len)

    return input_ids[0]

def encode_target_str(text, tokenizer, seq_len,
                    lang_token_map=LANG_TOKEN_MAPPING):
    token_ids = tokenizer.encode(
        text = text,
        return_tensors = 'pt',
        padding = 'max_length',
        truncation = True,
        max_length = seq_len)

    return token_ids[0]

def format_translation_data(translations, lang_token_map,
                          tokenizer, seq_len=128):
    # Choose a random 2 languages for in i/o
    langs = list(lang_token_map.keys())
    input_lang, target_lang = np.random.choice(langs, size=2, replace=False)

    # Get the translations for the batch
    input_text = translations[input_lang]
    target_text = translations[target_lang]

    if input_text is None or target_text is None:
        return None

    input_token_ids = encode_input_str(
        input_text, target_lang, tokenizer, seq_len, lang_token_map)

    target_token_ids = encode_target_str(
        target_text, tokenizer, seq_len, lang_token_map)

```

```

    return input_token_ids, target_token_ids

def transform_batch(batch, lang_token_map, tokenizer):
    inputs = []
    targets = []
    for translation_set in batch['translation']:
        formatted_data = format_translation_data(
            translation_set, lang_token_map, tokenizer, max_seq_len)

        if formatted_data is None:
            continue

        input_ids, target_ids = formatted_data
        inputs.append(input_ids.unsqueeze(0))
        targets.append(target_ids.unsqueeze(0))

    batch_input_ids = torch.cat(inputs).cuda()
    batch_target_ids = torch.cat(targets).cuda()

    return batch_input_ids, batch_target_ids

def get_data_generator(dataset, lang_token_map, tokenizer, batch_size=32):
    dataset = dataset.shuffle()
    for i in range(0, len(dataset), batch_size):
        raw_batch = dataset[i:i+batch_size]
        yield transform_batch(raw_batch, lang_token_map, tokenizer)

# Testing `data_transform`
in_ids, out_ids = format_translation_data(
    train_dataset[0]['translation'], LANG_TOKEN_MAPPING, tokenizer)

print(' '.join(tokenizer.convert_ids_to_tokens(in_ids)))
print(' '.join(tokenizer.convert_ids_to_tokens(out_ids)))

# Testing data generator
data_gen = get_data_generator(train_dataset, LANG_TOKEN_MAPPING, tokenizer, 8)
data_batch = next(data_gen)
print('Input shape:', data_batch[0].shape)
print('Output shape:', data_batch[1].shape)

"""# Train/Finetune BERT"""

model.load_state_dict(torch.load(model_path))

# Constants
n_epochs = 8
batch_size = 16

```

```

print_freq = 50
checkpoint_freq = 1000
lr = 5e-4
n_batches = int(np.ceil(len(train_dataset) / batch_size))
total_steps = n_epochs * n_batches
n_warmup_steps = int(total_steps * 0.01)

# Optimizer
optimizer = AdamW(model.parameters(), lr=lr)
scheduler = get_linear_schedule_with_warmup(
    optimizer, n_warmup_steps, total_steps)

losses = []

def eval_model(model, gdataset, max_iters=8):
    test_generator = get_data_generator(gdataset, LANG_TOKEN_MAPPING,
                                        tokenizer, batch_size)

    eval_losses = []
    for i, (input_batch, label_batch) in enumerate(test_generator):
        if i >= max_iters:
            break

        model_out = model.forward(
            input_ids = input_batch,
            labels = label_batch)
        eval_losses.append(model_out.loss.item())

    return np.mean(eval_losses)

for epoch_idx in range(n_epochs):
    # Randomize data order
    data_generator = get_data_generator(train_dataset, LANG_TOKEN_MAPPING,
                                        tokenizer, batch_size)

    for batch_idx, (input_batch, label_batch) \
        in tqdm_notebook(enumerate(data_generator), total=n_batches):
        optimizer.zero_grad()

        # Forward pass
        model_out = model.forward(
            input_ids = input_batch,
            labels = label_batch)

        # Calculate loss and update weights
        loss = model_out.loss
        losses.append(loss.item())
        loss.backward()
        optimizer.step()

```

```

scheduler.step()

# Print training update info
if (batch_idx + 1) % print_freq == 0:
    avg_loss = np.mean(losses[-print_freq:])
    print('Epoch: {} | Step: {} | Avg. loss: {:.3f} | lr: {}'.format(
        epoch_idx+1, batch_idx+1, avg_loss, scheduler.get_last_lr()[0]))

if (batch_idx + 1) % checkpoint_freq == 0:
    test_loss = eval_model(model, test_dataset)
    print('Saving model with test loss of {:.3f}'.format(test_loss))
    torch.save(model.state_dict(), model_path)

torch.save(model.state_dict(), model_path)

# Graph the loss

window_size = 50
smoothed_losses = []
for i in range(len(losses)-window_size):
    smoothed_losses.append(np.mean(losses[i:i+window_size]))

plt.plot(smoothed_losses[100:])

"""# Manual Testing"""

test_sentence = test_dataset[0]['translation']['en']
print('Raw input text:', test_sentence)

input_ids = encode_input_str(
    text = test_sentence,
    target_lang = 'ja',
    tokenizer = tokenizer,
    seq_len = model.config.max_length,
    lang_token_map = LANG_TOKEN_MAPPING)
input_ids = input_ids.unsqueeze(0).cuda()

print('Truncated input text:', tokenizer.convert_tokens_to_string(
    tokenizer.convert_ids_to_tokens(input_ids[0])))

output_tokens = model.generate(input_ids, num_beams=10,
    num_return_sequences=3)
# print(output_tokens)
for token_set in output_tokens:
    print(tokenizer.decode(token_set, skip_special_tokens=True))

#@title Slick Blue Translate

```

```

input_text = 'A surfboarder ran into a shark' #@param {type:"string"}
output_language = 'ja' #@param ["en", "ja", "zh"]

input_ids = encode_input_str(
    text = input_text,
    target_lang = output_language,
    tokenizer = tokenizer,
    seq_len = model.config.max_length,
    lang_token_map = LANG_TOKEN_MAPPING)
input_ids = input_ids.unsqueeze(0).cuda()

output_tokens = model.generate(input_ids, num_beams=20, length_penalty=0.2)
print(input_text + ' -> ' + \
      tokenizer.decode(output_tokens[0], skip_special_tokens=True))

```

Output:

input_text: " A surfboarder ran into a shark

output_language:

[Show code](#)

A surfboarder ran into a shark -> 水泳選手はサメの群れの中にいた。