

```
In [1]: import torch
import torch.nn as nn
import torch.nn.functional as F
from collections import Counter
from torch.utils.data import Dataset, DataLoader
import numpy as np
import re
device = 'cuda:0' if torch.cuda.is_available() else 'cpu'
```

```
In [2]: corpus = [line.strip() for line in open('TheTimeMachine.txt') if line.strip()]
print("\n".join(corpus[:10]))

# Tokenize the sentences into words
# All lower caps. Ignore punctuation.
corpus = [re.sub('[^A-Za-z0-9]+', ' ', line).lower() for line in corpus]
corpus = [re.sub(' +', ' ', line) for line in corpus]
corpus = [word for line in corpus for word in line.split()]
```

The Time Machine, by H. G. Wells [1898]

I

The Time Traveller (for so it will be convenient to speak of him) was expounding a recondite matter to us. His grey eyes shone and twinkled, and his usually pale face was flushed and animated. The fire burned brightly, and the soft radiance of the incandescent lights in the lilies of silver caught the bubbles that flashed and passed in our glasses. Our chairs, being his patents, embraced and caressed us rather than submitted to be sat upon, and there was that luxurious after-dinner atmosphere when thought roams gracefully

```
In [3]: vocab_size = 5000
tkn_counter = Counter([word for word in corpus])
vocab = {word: idx for idx, (word, _) in enumerate(tkn_counter.most_common(vocab_size))}
vocab["/UNK"] = len(vocab)
print(f" * Found {len(vocab)} unique words in the provided corpus (of size {vocab_size})")
print(f" * Created vocabulary from corpus.\n")
print(f" * The 10 most common words are the following:")
print(tkn_counter.most_common(10))
```

* Found 4582 unique words in the provided corpus (of size 32776).

* Created vocabulary from corpus.

* The 10 most common words are the following:

```
[('the', 2261), ('i', 1267), ('and', 1245), ('of', 1155), ('a', 816), ('to', 695), ('was', 552), ('in', 541), ('that', 443), ('my', 440)]
```

```
In [7]: class TextCorpusDataset(Dataset):
    def __init__(self, corpus, vocab, sequence_len=50):
        super().__init__()
        self.corpus = corpus
        self.sequence_len = sequence_len

    # Vocabulary (word-to-index mapping)
    self.vocab = vocab
```

```

# Inverse vocabulary (index-to-word mapping)
self.inv_vocab = {idx: word for word, idx in self.vocab.items()}

def convert2idx(self, word_sequence):
    return [self.vocab[word] if word in self.vocab else "/UNK" for word in word_sequence]

def convert2words(self, idx_sequence):
    return [self.inv_vocab[idx] for idx in idx_sequence]

def __len__(self):
    return (len(self.corpus) - self.sequence_len) // self.sequence_len

def __getitem__(self, idx):
    idx = idx * self.sequence_len
    snippet = self.corpus[idx:idx+self.sequence_len]
    snippet = torch.tensor(self.convert2idx(snippet))
    return snippet

# Test dataset function
dataset = TextCorpusDataset(corpus, vocab, sequence_len=50)
sequence = dataset[4]
print(len(dataset))
print("\nRandom sequence from the corpus.")
print(" * Token IDS:\t", sequence)
print(" * Words:\t\t", " ".join([dataset.inv_vocab[i] for i in sequence.tolist()]))

```

654

Random sequence from the corpus.

```

 * Token IDS:  tensor([ 21,    5, 2211,  682,  275, 1430,  235,   15,   1
0,   21,  114,  196,
                830,   13,  180,   13,    1,  502,   29,   21,   21,  150,    3,   3
12,
                8,    4, 1047,  330,    4,  330,    3,  683, 2212,  187,   42,   4
00,
                591,   28, 1427,   21,    8, 1431,  187,    4, 1047, 2213,   58,   1
32,
                90,  244])

```

* Words: you to accept anything without reasonable ground for it you will soon admit as much as i need from you you know of course that a mathematical line a line of thickness nil has no real existence they taught you that neither has a mathematical plane these things are mere

```

In [ ]: class SimpleRNN(nn.Module):
        """A RNN Model implemented from scratch."""
        def __init__(self, vocab_size, hidden_dim):
            super().__init__()
            self.vocab_size, self.hidden_dim = vocab_size, hidden_dim

            self.inp2state = nn.Linear(vocab_size, hidden_dim)
            self.state2state = nn.Linear(hidden_dim, hidden_dim)
            self.state2out = nn.Linear(hidden_dim, vocab_size)

```

```

        for m in self.modules():
            if isinstance(m, nn.Linear):
                nn.init.normal_(m.weight, std=0.01)
                nn.init.zeros_(m.bias)

    def initial_state(self, batch_size, device):
        return torch.zeros((batch_size, self.hidden_dim)).to(device)

    def forward(self, inp_seq, state=None):
        n_steps, batch_size = inp_seq.shape[:2]

        # If state is not provided, get initial state.
        if state is None:
            state = self.initial_state(batch_size, inp_seq.device)

        outputs = []
        for t in range(n_steps):
            inp_at_t = inp_seq[t]
            # Compute new state:  $h_t = \tanh(W_{x2h} * x_t + W_{h2h} * h_{t-1} + b_h)$ 
            state = torch.tanh(self.inp2state(inp_at_t) + self.state2state(s
            # Compute output:  $o_t = W_{h2o} * h_t + b_o$ 
            out = self.state2out(state)
            outputs.append(out)
        outputs = torch.stack(outputs, 0)

        return outputs, state

hidden_dim = 256
model = SimpleRNN(len(vocab), hidden_dim).to(device)

```

```

In [30]: sentence = "today is too darn cold".split()
inp = F.one_hot(torch.tensor(dataset.convert2idx(sentence), device=device),
inp = inp.unsqueeze(1)
print(inp.shape)
Yhat, _ = model(inp)
print(Yhat.shape)
predicted_indices = torch.argmax(Yhat, dim=-1).squeeze(1)
Yhat_words = dataset.convert2words(predicted_indices.tolist())
print(Yhat_words)

torch.Size([5, 1, 4582])
torch.Size([5, 1, 4582])
['fix', 'nature', 'collapsed', 'fix', 'endowed']

```

```

In [29]: @torch.no_grad()
def generate(prefix, num_preds, model, vocab):
    """Generates a sentence following the `prefix`."""
    prefix = torch.tensor(dataset.convert2idx(prefix.split()), device=device)

    state, outputs = None, [prefix[0]]
    for i in range(1, len(prefix) + num_preds):
        # Prepare the current token to feed the model
        inp = F.one_hot(outputs[-1], len(vocab)).float()
        inp = inp[None, None]

```

```

# Compute the prediction of the next token
yhat, state = model(inp, state)

if i < len(prefix):
    # During warmup (while parsing the prefix), we ignore the model
    outputs.append(prefix[i])
else:
    # Otherwise, append the model prediction to the list
    yhat = yhat[..., :-1].argmax(dim=-1).reshape(1).long()
    outputs.append(yhat)
return ' '.join([dataset.inv_vocab[tkn.item()] for tkn in outputs])

generate('i do not mean to ask you to accept anything', 10, model, vocab)

```

Out[29]: 'i do not mean to ask you to accept anything tampering palaeontology expression declaration concerned hastings servants hull explain impartiality'

```

In [ ]: def train_on_sequence(seq, model, optimizer, unroll=5):
    """Train the model within a batch of long text sequences."""
    batch_size, num_tokens = seq.shape

    total_loss, state = 0., None
    for i in range(0, num_tokens-unroll-1, unroll):
        if state is not None:
            state.detach_()

        # Define the input sequence along which we will unroll the RNN
        x_inp = seq[:, i:i+unroll] # Must be of size T x B
        y_trg = seq[:, i+1:i+unroll+1] # Must be of size T x B

        # Forward the model and compute the loss
        x_inp = F.one_hot(x_inp, len(vocab)).float()
        y_hat, state = model(x_inp, state)
        l = loss(y_hat.flatten(0, 1), y_trg.flatten(0, 1).long())
        total_loss += l.item()

        # Backward step
        optimizer.zero_grad()
        l.backward()
        optimizer.step()

    n_batches = (num_tokens-unroll-1) // unroll
    return total_loss/n_batches

def fit(model, loader, vocab, lr, num_epochs=100, unroll=5):
    optimizer = torch.optim.SGD(model.parameters(), lr, momentum=0.9)
    test_prompt = 'i do not mean to ask you to accept anything'
    for epoch in range(num_epochs):
        total_loss = 0
        for sequence in loader:
            total_loss += train_on_sequence(sequence.to(device), model, optimizer)
        total_loss /= len(loader)

```


