

Recurrent Neural Networks (RNN)

© 李浩东 3190104890@zju.edu.cn

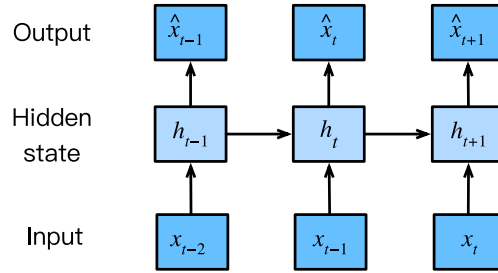
- Sequence Models
- Text Preprocessing
- Language Models and the Dataset
- Recurrent Neural Networks

$$x_t \sim P(x_t | x_{t-1}, \dots, x_1) \quad (1)$$

- Music, speech, text, and videos are all sequential in nature. If we were to permute them they would make little sense. The headline dog bites man is much less surprising than man bites dog, even though the words are identical.
- In order to achieve this, we could use a regression model.

$$x_t \sim P(x_t | x_{t-1}, \dots, x_1) \sim P(x_t | x_{t-1}, x_{t-2}, \dots, x_{t-\tau}) \quad (2)$$

- Such models will be called autoregressive models, as they quite literally perform regression on themselves.
- The second strategy, shown in the figure below, is to keep some summary h_t of the past observations, and at the same time update h_t in addition to the prediction \hat{x}_t . This leads to models that estimate x_t with $\hat{x}_t = P(x_t | h_t)$ and moreover updates of the form $h_t = g(h_{t-1}, x_{t-1})$. Since h_t is never observed, these models are also called latent autoregressive models.



- Both cases raise the obvious question of how to generate training data. One typically uses historical observations to predict the next observation given the ones up to right now. Obviously we do not expect time to stand still. However, a common assumption is that while the specific values of x_t might change, at least the dynamics of the sequence itself will not. This is reasonable, since novel dynamics are just that, novel and thus not predictable using data that we have so far. Statisticians call dynamics that do not change stationary. Regardless of what we do, we will thus get an estimate of the entire sequence via:

$$P(x_1, \dots, x_T) = \prod_{t=1}^T P(x_t | x_{t-1}, \dots, x_1) \quad (3)$$

- Note that the above considerations still hold if we deal with discrete objects, such as words, rather than continuous numbers. The only difference is that in such a situation we need to use a classifier rather than a regression model to estimate it.

Markov Models

- Recall the approximation that in an autoregressive model we use only $x_{t-1}, \dots, x_{t-\tau}$ instead of x_{t-1}, \dots, x_{t-1} to estimate x_t . Whenever this approximation is accurate we say that the sequence satisfies a Markov condition. In particular, if $\tau = 1$, we have a first-order Markov model and $P(x)$ is given by:

$$P(x_1, \dots, x_T) = \prod_{t=1}^T P(x_t | x_{t-1}) \text{ where } P(x_1 | x_0) = P(x_1) \quad (4)$$

- Such models are particularly nice whenever x_t assumes only a discrete value, since in this case dynamic programming can be used to compute values along the chain exactly. For instance, we can compute $P(x_{t+1} | x_{t-1})$ efficiently:

$$\begin{aligned} P(x_{t+1} | x_{t-1}) &= \frac{\sum_{x_t} P(x_{t+1}, x_t, x_{t-1})}{P(x_{t-1})} \\ &= \frac{\sum_{x_t} P(x_{t+1} | x_t, x_{t-1}) P(x_t, x_{t-1})}{P(x_{t-1})} \\ &= \sum P(x_{t+1} | x_t) P(x_t | x_{t-1}) \end{aligned} \quad (5)$$

x_t

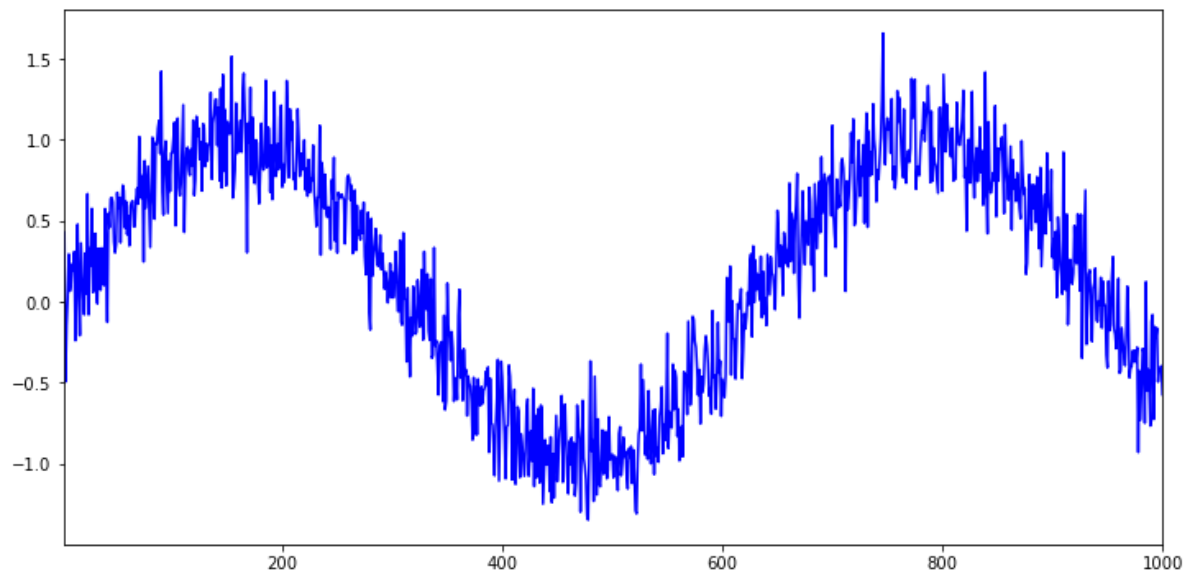
- In fact, if we have a Markov model, we can obtain a reverse conditional probability distribution, too. In many cases, however, there exists a natural direction for the data, namely going forward in time. It is clear that future events cannot influence the past. Hence, if we change x_t , we may be able to influence what happens for x_{t+1} going forward but not the converse. That is, if we change x_t , the distribution over past events will not change. Consequently, it ought to be easier to explain $P(x_{t+1} | x_t)$ rather than $P(x_t | x_{t+1})$.

```
1 import torch
2 from torch import nn
3 import matplotlib.pyplot as plt
4 from torch.utils.data import DataLoader
5 import torch.backends.cudnn as cudnn
6 import torch.optim as optim
7 import requests
8 import collections
9 from torch.nn import functional as F
10 import re
11 # import time
12 import math
13 import random
14 from torch.utils import data
15 import hashlib
16 import os
17 print(torch.__version__)
```

```
1 1.11.0
```

```
1 T = 1000 # Generate a total of 1000 points
2 time = torch.arange(1, T + 1, dtype=torch.float32)
3 x = torch.sin(0.01 * time) + torch.normal(0, 0.2, (T,))
4 plt.figure(figsize=(12, 6))
5 plt.plot(time, x, '-b')
6 plt.xlim(1, 1000)
```

```
1 (1.0, 1000.0)
```



```

1 def generate_date(tau):
2     features = torch.zeros((T - tau, tau))
3     for i in range(tau):
4         features[:, i] = x[i: T - tau + i]
5     labels = x[tau:].reshape((-1, 1))
6     num_workers = 0
7     batch_size, n_train = 20, 600
8     train_loader = DataLoader(dataset=data.TensorDataset(features[:n_train], labels[:n_train]),
9                             batch_size=batch_size, shuffle=True, num_workers=num_workers)
10    train_iter = iter(train_loader)
11    test_loader = DataLoader(dataset=data.TensorDataset(features[n_train:], labels[n_train:]),
12                            batch_size=batch_size, shuffle=True, num_workers=num_workers)
13    test_iter = iter(test_loader)
14    print(len(train_loader), len(test_loader))
15    print(len(train_iter), len(test_iter))
16    return train_loader, test_loader, features
17
18 train_loader, test_loader, features = generate_date(tau=4)

```

```

1 30 20
2 30 20

```

```

1 datas, targets = next(iter(train_loader))
2 print(datas)
3 print(targets)

```

```

1 tensor([[ 0.9129,  1.4217,  0.6786,  0.5357],
2         [ 0.3694,  0.3018,  0.3527,  0.6734],
3         [ 0.1396, -0.3484, -0.3318,  0.3335],
4         [ 1.2910,  0.7567,  1.1020,  1.1510],
5         [-0.6803, -0.6971, -0.7796, -0.3869],
6         [-0.6435, -0.5256, -0.6300, -0.5475],
7         [ 0.2834,  0.3532,  0.4497,  0.2817],
8         [ 0.9646,  1.1996,  0.7468,  1.3165],
9         [-1.0178, -1.0986, -1.3208, -1.3482],
10        [ 0.1829,  0.1130, -0.0066,  0.0979],
11        [ 0.0174, -0.2091,  0.3615,  0.1878],
12        [-0.8771, -1.0745, -1.1865, -0.8613],
13        [-0.6647, -0.9520, -0.8991, -0.9913],
14        [ 1.1012,  0.9812,  1.3224,  0.7805],
15        [ 0.7449,  0.9660,  0.6778,  0.5861],
16        [ 0.7645,  0.6381,  0.7783,  0.2462],
17        [ 0.3315,  0.2322,  0.1001,  0.5734],
18        [ 0.6303,  0.3578,  0.5908,  0.7649],
19        [ 0.9270,  0.9531,  0.8100,  0.7802],
20        [-1.0004, -0.8097, -0.6716, -1.0671]])
21 tensor([ 0.9197,
22         0.5742,
23        -0.2718,
24         1.2469,
25        -0.6670,
26        -0.5339,
27         0.2200,
28         0.7036,
29        -0.9049,
30         0.2362,
31         0.0244,
32        -1.0070,
33        -0.7042,
34         1.1377,
35         0.4639],

```

```

36         [ 0.8699],
37         [ 0.2883],
38         [ 0.7811],
39         [ 0.9399],
40         [-0.6647]])

```

```

1  def simple_mlp():
2      net = nn.Sequential(nn.Linear(4, 32),
3                          nn.ReLU(),
4                          nn.Linear(32, 1))
5      return net
6
7  class SimpleMLP(nn.Module):
8      def __init__(self):
9          super(SimpleMLP, self).__init__()
10         self.model = simple_mlp()
11     def forward(self, x):
12         x = self.model(x)
13         return x

```

```

1  def train_model(model, name, epoch_num=500, device='cuda', learning_rate=0.01,
2                  train_loader=train_loader, test_loader=test_loader):
3      def init_weights(m):
4          if type(m) == nn.Linear or type(m) == nn.Conv2d:
5              nn.init.xavier_uniform_(m.weight)
6      model.apply(init_weights)
7      if device == 'cuda':
8          print("Use CUDA for training.")
9          model = torch.nn.DataParallel(model) # make parallel
10         cudnn.benchmark = True
11
12     model.to(device)
13     # specify loss function
14     criterion = nn.MSELoss(reduction='none')
15     # specify optimizer
16     optimizer = optim.Adam(model.parameters(), learning_rate)
17     train_losslist = []
18     valid_losslist = []
19
20     for epoch in range(1, epoch_num+1):
21         # keep track of training and validation loss
22         train_loss = 0.0
23         valid_loss = 0.0
24
25         model.train()
26         for data, target in train_loader:
27             data, target = data.to(device), target.to(device)
28             optimizer.zero_grad()
29             output = model(data)
30             loss = criterion(output, target)
31             loss.sum().backward()
32             optimizer.step()
33             train_loss += loss.sum().item()*data.size(0)
34
35         model.eval()
36         for data, target in test_loader:
37             data, target = data.to(device), target.to(device)
38             output = model(data)
39             loss = criterion(output, target)
40             valid_loss += loss.sum().item()*data.size(0)
41
42     train_loss = train_loss/len(train_loader.dataset)
43     valid_loss = valid_loss/len(test_loader.dataset)
44     train_losslist.append(train_loss)
45     valid_losslist.append(valid_loss)

```

```

45     print("Epoch ->", epoch, "\t train_loss ->", train_loss, "\t\t\t valid_loss ->",
valid_loss)
46
47     os.mkdir("./result/") if os.path.exists("./result/") == False else print("./result/ exists.")
48     torch.save(model.state_dict(), './result/' + name + '.pt')
49     return train_losslist, valid_losslist, model

```

```

1 device = "cuda" if torch.cuda.is_available() else "cpu"
2 print(torch.cuda.get_arch_list(), device)
3 train_losslist, valid_losslist, model = train_model(SimpleMLP(), "reg_sin_simplemlp",
4                                                     10, device, learning_rate=0.05)

```

```

1 ['sm_37', 'sm_50', 'sm_60', 'sm_61', 'sm_70', 'sm_75', 'sm_80', 'sm_86', 'compute_37'] cuda
2 Use CUDA for training.
3 Epoch -> 1   train_loss -> 2.590431062380473           valid_loss -> 1.2389005639336326
4 Epoch -> 2   train_loss -> 1.1704057296117147         valid_loss -> 1.7951183860952205
5 Epoch -> 3   train_loss -> 1.2141841918230056         valid_loss -> 1.2576977171079078
6 Epoch -> 4   train_loss -> 1.1336304724216462         valid_loss -> 1.1760255363252428
7 Epoch -> 5   train_loss -> 1.1412086109320323         valid_loss -> 1.1651216161371483
8 Epoch -> 6   train_loss -> 1.2022289117177327         valid_loss -> 1.1709889817719508
9 Epoch -> 7   train_loss -> 1.1765139768520991         valid_loss -> 1.1640310149000148
10 Epoch -> 8   train_loss -> 1.2275258978207906         valid_loss -> 1.397282317431286
11 Epoch -> 9   train_loss -> 1.3512031773726145         valid_loss -> 1.2245428495936923
12 Epoch -> 10   train_loss -> 1.1404260714848837         valid_loss -> 1.2579764302330787
13 ./result/ exists.

```

```

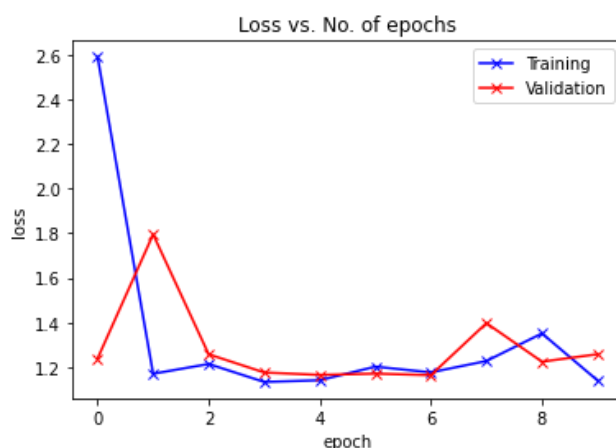
1 def plot_losses(train_losslist, valid_losslist):
2     plt.plot(train_losslist, '-bx')
3     plt.plot(valid_losslist, '-rx')
4     plt.xlabel('epoch')
5     plt.ylabel('loss')
6     plt.legend(['Training', 'Validation'])
7     plt.title('Loss vs. No. of epochs')
8     plt.show()

```

```

1 plot_losses(train_losslist, valid_losslist)

```

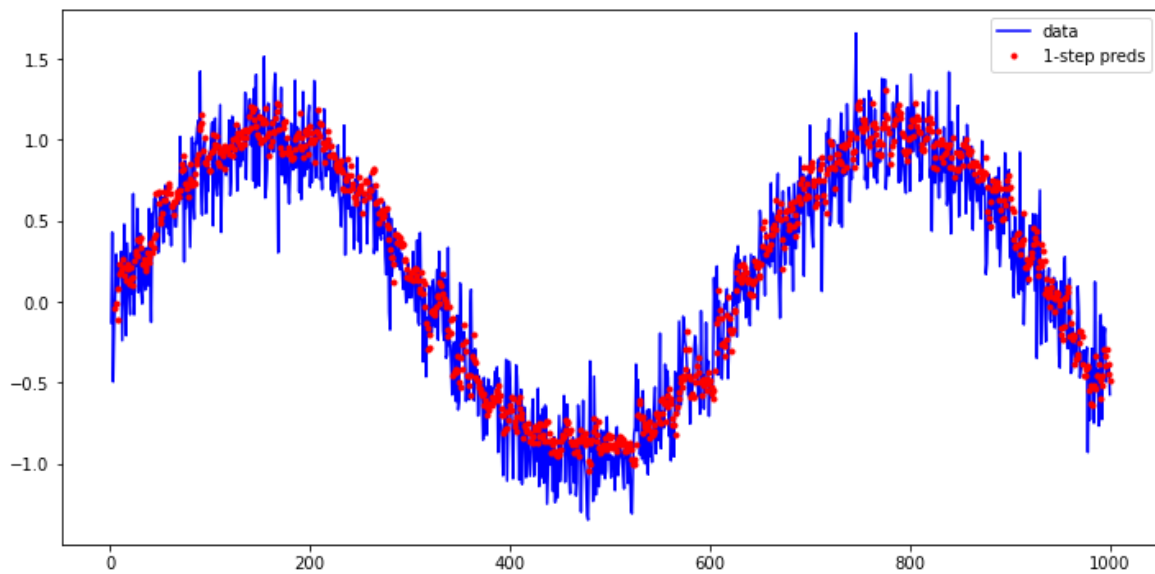


```

1 tau, n_train = 4, 600
2 onestep_preds = model(features)
3 plt.figure(figsize=(12, 6))
4 plt.plot(time, x.detach().numpy(), '-b')
5 plt.plot(time[tau:], onestep_preds.cpu().detach().numpy(), '.r')
6 plt.legend(['data', '1-step preds'])

```

```
1 | <matplotlib.legend.Legend at 0x1f7017ba5b0>
```



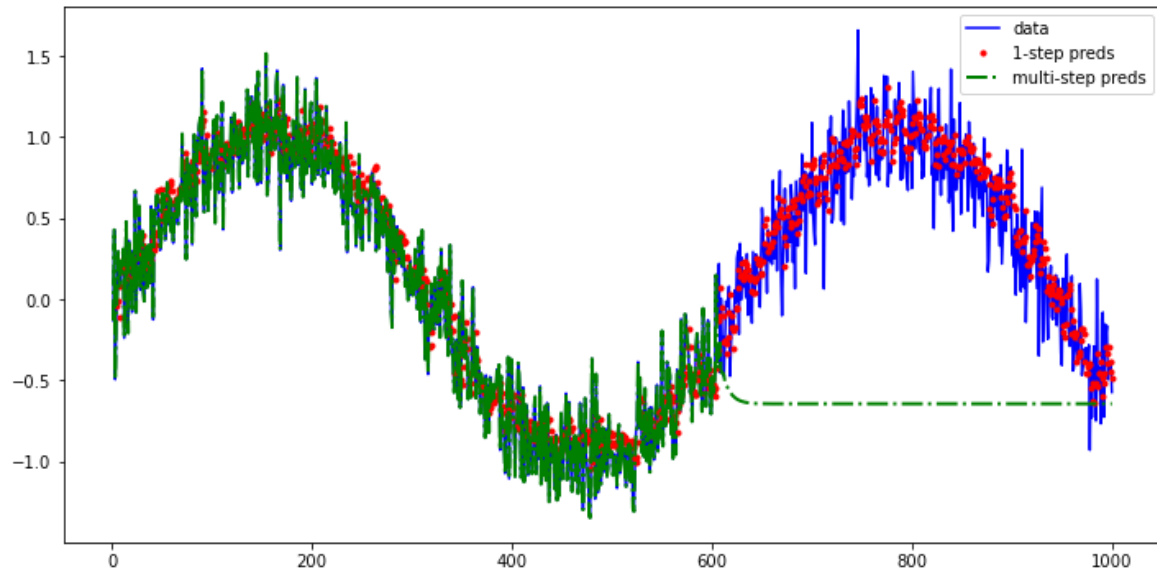
- The one-step-ahead predictions look nice, just as we expected. Even beyond 604 (`n_train + tau`) observations the predictions still look trustworthy. However, there is just one little problem to this: if we observe sequence data only until time step 604, we cannot hope to receive the inputs for all the future one-step-ahead predictions. Instead, we need to work our way forward one step at a time:

$$\begin{aligned}
 \hat{x}_{605} &= f(x_{601}, x_{602}, x_{603}, x_{604}), \\
 \hat{x}_{606} &= f(x_{602}, x_{603}, x_{604}, \hat{x}_{605}), \\
 \hat{x}_{607} &= f(x_{603}, x_{604}, \hat{x}_{605}, \hat{x}_{606}), \\
 \hat{x}_{608} &= f(x_{604}, \hat{x}_{605}, \hat{x}_{606}, \hat{x}_{607}), \\
 \hat{x}_{609} &= f(\hat{x}_{605}, \hat{x}_{606}, \hat{x}_{607}, \hat{x}_{608}), \\
 &\dots
 \end{aligned}
 \tag{6}$$

```
1 | multistep_preds = torch.zeros(T)
2 | multistep_preds[: n_train + tau] = x[: n_train + tau]
3 | for i in range(n_train + tau, T):
4 |     multistep_preds[i] = model(multistep_preds[i - tau:i].reshape((1, -1)))
```

```
1 | plt.figure(figsize=(12, 6))
2 | plt.plot(time, x.detach().numpy(), '-b')
3 | plt.plot(time[tau:], onestep_preds.cpu().detach().numpy(), '.r')
4 | plt.plot(time, multistep_preds.cpu().detach().numpy(), '-.g', linewidth=2)
5 | plt.legend(['data', '1-step preds', 'multi-step preds'])
6 | # plt.ylim(-4, 2)
```

```
1 | <matplotlib.legend.Legend at 0x1f701805550>
```



- Let us take a closer look at the difficulties in k -step-ahead predictions by computing predictions on the entire sequence for $k = 1, 4, 16, 64$.

```

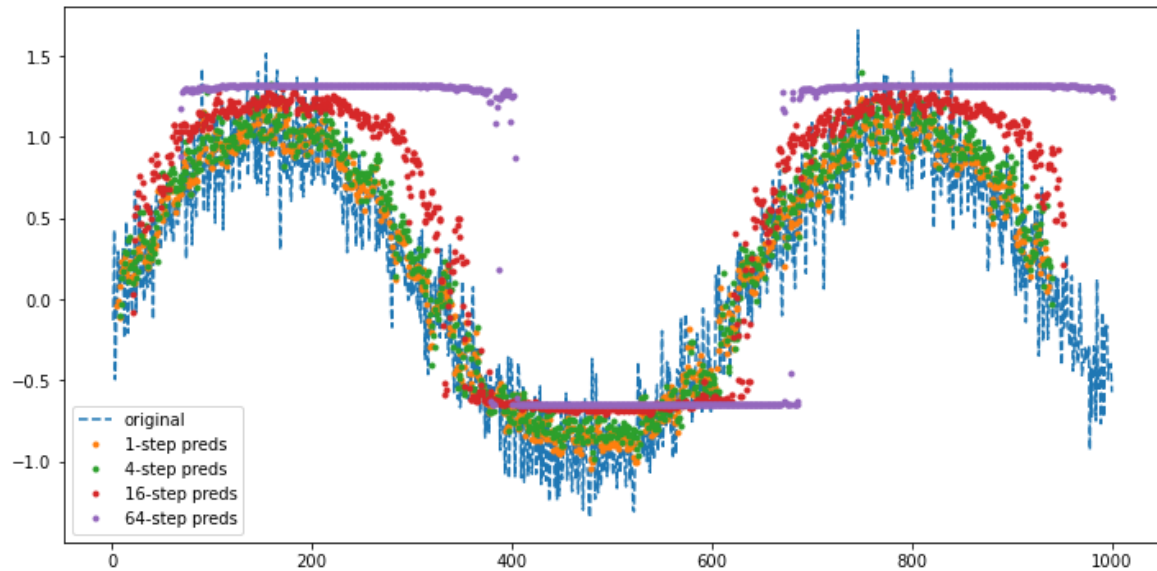
1 max_steps = 64
2
3 features = torch.zeros((T - tau - max_steps + 1, tau + max_steps))
4 # Column `i` (`i` < `tau`) are observations from `x` for time steps from
5 # `i + 1` to `i + T - tau - max_steps + 1`
6 for i in range(tau):
7     features[:, i] = x[i: i + T - tau - max_steps + 1]
8
9 # Column `i` (`i` >= `tau`) are the `(i - tau + 1)`-step-ahead predictions for
10 # time steps from `i + 1` to `i + T - tau - max_steps + 1`
11 for i in range(tau, tau + max_steps):
12     features[:, i] = model(features[:, i - tau:i]).reshape(-1)
13
14 steps = (1, 4, 16, 64)
15 plt.figure(figsize=(12, 6))
16 plt.plot(time, x.detach().numpy(), '--')
17 for step in steps:
18     plt.plot(time[tau + step - 1: T - max_steps + step], features[:, (tau + step -
19     1)].cpu().detach().numpy(), '.')
19 plt.legend(["original"] + [str(step) + '-step preds' for step in steps])

```

```

1 <matplotlib.legend.Legend at 0x1f704cd7f40>

```



Text Preprocessing

- Load text as strings into memory.
- Split strings into tokens (e.g., words and characters).
- Build a table of vocabulary to map the split tokens to numerical indices.
- Convert text into sequences of numerical indices so they can be manipulated by models easily.

```

1 def download(url, sha1_hash, cache_dir=os.path.join('.', 'data')):
2     os.makedirs(cache_dir, exist_ok=True)
3     fname = os.path.join(cache_dir, url.split('/')[-1])
4     if os.path.exists(fname):
5         sha1 = hashlib.sha1()
6         with open(fname, 'rb') as f:
7             while True:
8                 data = f.read(1048576)
9                 if not data:
10                    break
11                sha1.update(data)
12            if sha1.hexdigest() == sha1_hash:
13                return fname # Hit cache
14    print(f'Downloading {fname} from {url}...')
15    r = requests.get(url, stream=True, verify=True)
16    with open(fname, 'wb') as f:
17        f.write(r.content)
18    return fname

```

```

1 txt_file = download('http://d21-data.s3-accelerate.amazonaws.com/' + 'timemachine.txt',
2                     '090b5e7e70c295757f55df93cb0a180b9691891a')
3 def read_time_machine(): #@save
4     """Load the time machine dataset into a list of text lines."""
5     with open(txt_file, 'r') as f:
6         lines = f.readlines()
7         return [re.sub('[A-Za-z]+', ' ', line).strip().lower() for line in lines]
8
9 lines = read_time_machine()
10 print("text lines ->", len(lines))
11 print(lines[0])
12 print(lines[10])

```



```

1 | text lines -> 3221
2 | the time machine by h g wells
3 | twinkled and his usually pale face was flushed and animated the

```

Tokenization

- The following tokenize function takes a list (lines) as the input, where each element is a text sequence (e.g., a text line). Each text sequence is split into a list of tokens. A token is the basic unit in text. In the end, a list of token lists are returned, where each token is a string.

```

1 | def tokenize(lines, token='char'): #@save
2 |     """Split text lines into word or character tokens."""
3 |     if token == 'word':
4 |         return [line.split() for line in lines]
5 |     elif token == 'char':
6 |         return [list(line) for line in lines]
7 |     else:
8 |         print('ERROR: unknown token type: ' + token)
9 |
10 | tokens = tokenize(lines, token='word')
11 | for i in range(11):
12 |     print(tokens[i])

```

```

1 | ['the', 'time', 'machine', 'by', 'h', 'g', 'wells']
2 | []
3 | []
4 | []
5 | []
6 | ['i']
7 | []
8 | []
9 | ['the', 'time', 'traveller', 'for', 'so', 'it', 'will', 'be', 'convenient', 'to', 'speak', 'of',
10 |  'him']
11 | ['was', 'expounding', 'a', 'recondite', 'matter', 'to', 'us', 'his', 'grey', 'eyes', 'shone',
    'and']
12 | ['twinkled', 'and', 'his', 'usually', 'pale', 'face', 'was', 'flushed', 'and', 'animated', 'the']

```

Vocabulary

- The string type of the token is inconvenient to be used by models, which take numerical inputs. Now let us build a dictionary, often called vocabulary as well, to map string tokens into numerical indices starting from 0. To do so, we first count the unique tokens in all the documents from the training set, namely a corpus, and then assign a numerical index to each unique token according to its frequency. Rarely appeared tokens are often removed to reduce the complexity. Any token that does not exist in the corpus or has been removed is mapped into a special unknown token "<unk>". We optionally add a list of reserved tokens, such as "<pad>" for padding, "<bos>" to present the beginning for a sequence, and "<eos>" for the end of a sequence.

```

1 | class Vocab: #@save
2 |     """Vocabulary for text."""
3 |     def __init__(self, tokens=None, min_freq=0, reserved_tokens=None):
4 |         if tokens is None:
5 |             tokens = []
6 |         if reserved_tokens is None:
7 |             reserved_tokens = []
8 |         # Sort according to frequencies
9 |         counter = count_corpus(tokens)
10 |         self._token_freqs = sorted(counter.items(), key=lambda x: x[1],
11 |                                     reverse=True)
12 |         # The index for the unknown token is 0
13 |         self.idx_to_token = ['<unk>'] + reserved_tokens
14 |         self.token_to_idx = {token: idx
15 |                               for idx, token in enumerate(self.idx_to_token)}
16 |         for token, freq in self._token_freqs:
17 |             if freq < min_freq:

```

```

18         break
19         if token not in self.token_to_idx:
20             self.idx_to_token.append(token)
21             self.token_to_idx[token] = len(self.idx_to_token) - 1
22
23     def __len__(self):
24         return len(self.idx_to_token)
25
26     def __getitem__(self, tokens):
27         if not isinstance(tokens, (list, tuple)):
28             return self.token_to_idx.get(tokens, self.unk)
29         return [self.__getitem__(token) for token in tokens]
30
31     def to_tokens(self, indices):
32         if not isinstance(indices, (list, tuple)):
33             return self.idx_to_token[indices]
34         return [self.idx_to_token[index] for index in indices]
35
36     @property
37     def unk(self): # Index for the unknown token
38         return 0
39
40     @property
41     def token_freqs(self): # Index for the unknown token
42         return self._token_freqs
43
44     def count_corpus(tokens): #@save
45         """Count token frequencies."""
46         # Here `tokens` is a 1D list or 2D list
47         if len(tokens) == 0 or isinstance(tokens[0], list):
48             # Flatten a list of token lists into a list of tokens
49             tokens = [token for line in tokens for token in line]
50         return collections.Counter(tokens)

```

```

1 vocab = Vocab(tokens)
2 print(list(vocab.token_to_idx.items())[:10])
3 for i in [0, 10]:
4     print('words:', tokens[i])
5     print('indices:', vocab[tokens[i]])

```

```

1 [(<unk>, 0), ('the', 1), ('i', 2), ('and', 3), ('of', 4), ('a', 5), ('to', 6), ('was', 7), ('in',
2 8), ('that', 9)]
3 words: ['the', 'time', 'machine', 'by', 'h', 'g', 'wells']
4 indices: [1, 19, 50, 40, 2183, 2184, 400]
5 words: ['twinkled', 'and', 'his', 'usually', 'pale', 'face', 'was', 'flushed', 'and', 'animated',
6 'the']
7 indices: [2186, 3, 25, 1044, 362, 113, 7, 1421, 3, 1045, 1]

```

```

1 def load_corpus_time_machine(max_tokens=-1, token='char'): #@save
2     """Return token indices and the vocabulary of the time machine dataset."""
3     lines = read_time_machine()
4     tokens = tokenize(lines, token)
5     vocab = Vocab(tokens)
6     # Since each text line in the time machine dataset is not necessarily a
7     # sentence or a paragraph, flatten all the text lines into a single list
8     corpus = [vocab[token] for line in tokens for token in line]
9     if max_tokens > 0:
10         corpus = corpus[:max_tokens]
11     return corpus, vocab
12
13 corpus, vocab = load_corpus_time_machine(token='char')
14 print(len(corpus), len(vocab))
15 corpus, vocab = load_corpus_time_machine(token='word')
16 print(len(corpus), len(vocab))

```

Language Models

- Assume that the tokens in a text sequence of length T are in turn x_1, x_2, \dots, x_T . Then, in the text sequence, $x_t (1 \leq t \leq T)$ can be considered as the observation or label at time step t . Given such a text sequence, the goal of a language model is to estimate the joint probability of the sequence:

$$P(x_1, x_2, \dots, x_T) \quad (7)$$

- The obvious question is how we should model a document, or even a sequence of tokens. Suppose that we tokenize text data at the word level. Let us start by applying basic probability rules:

$$P(x_1, x_2, \dots, x_T) = \prod_{t=1}^T P(x_t | x_1, \dots, x_{t-1}) \quad (8)$$

$$P(\text{deep, learning, is, fun}) = P(\text{deep})P(\text{learning} | \text{deep})P(\text{is} | \text{deep, learning})P(\text{fun} | \text{deep, learning, is}) \quad (9)$$

- The probability of words can be calculated from the relative word frequency of a given word in the training dataset. For example, the estimate $\hat{P}(\text{deep})$ can be calculated as the probability of any sentence starting with the word “deep”. A slightly less accurate approach would be to count all occurrences of the word “deep” and divide it by the total number of words in the corpus.

$$\hat{P}(\text{learning} | \text{deep}) = \frac{n(\text{deep, learning})}{n(\text{deep})} \quad (10)$$

- where $n(x)$ and $n(x, x')$ are the number of occurrences of singletons and consecutive word pairs, respectively.
- Unfortunately, estimating the probability of a word pair is somewhat more difficult, since the occurrences of “deep learning” are a lot less frequent. In particular, for some unusual word combinations it may be tricky to find enough occurrences to get accurate estimates.
- Things take a turn for the worse for three-word combinations and beyond. There will be many plausible three-word combinations that we likely will not see in our dataset. Unless we provide some solution to assign such word combinations nonzero count, we will not be able to use them in a language model. If the dataset is small or if the words are very rare, we might not find even a single one of them.
- A common strategy is to perform some form of Laplace smoothing. The solution is to add a small constant to all counts. Denote by n the total number of words in the training set and m the number of unique words. This solution helps with singletons, e.g., via:

$$\begin{aligned} \hat{P}(x) &= \frac{n(x) + \epsilon_1/m}{n + \epsilon_1} \\ \hat{P}(x' | x) &= \frac{n(x, x') + \epsilon_2 \hat{P}(x')}{n(x) + \epsilon_2} \\ \hat{P}(x'' | x, x') &= \frac{n(x, x', x'') + \epsilon_3 \hat{P}(x'')}{n(x, x') + \epsilon_3} \end{aligned} \quad (11)$$

- Here ϵ_1, ϵ_2 and ϵ_3 are hyperparameters.
- Unfortunately, models like this get unwieldy rather quickly for the following reasons. First, we need to store all counts. Second, this entirely ignores the meaning of the words.
- For instance, “cat” and “feline” should occur in related contexts. It is quite difficult to adjust such models to additional contexts, whereas, deep learning based language models are well suited to take this into account.
- Last, long word sequences are almost certain to be novel, hence a model that simply counts the frequency of previously seen word sequences is bound to perform poorly there.

Markov Models

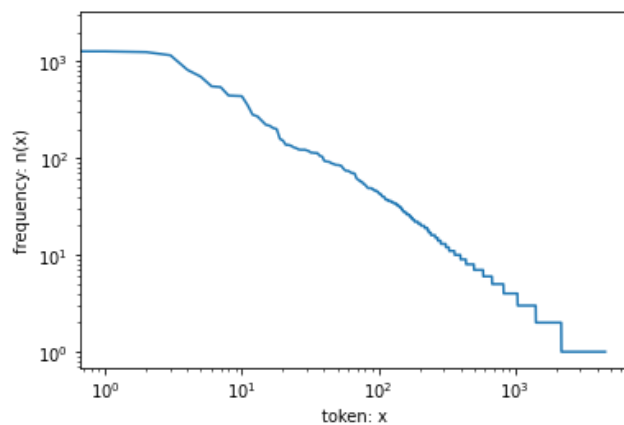
$$\begin{aligned} P(x_1, x_2, x_3, x_4) &= P(x_1)P(x_2)P(x_3)P(x_4) \\ P(x_1, x_2, x_3, x_4) &= P(x_1)P(x_2 | x_1)P(x_3 | x_2)P(x_4 | x_3) \\ P(x_1, x_2, x_3, x_4) &= P(x_1)P(x_2 | x_1)P(x_3 | x_1, x_2)P(x_4 | x_2, x_3) \end{aligned} \quad (12)$$

- The probability formulae that involve one, two, and three variables are typically referred to as unigram, bigram, and trigram models, respectively.

```
1 | vocab.token_freqs[:10]
```

```
1 | [('the', 2261),  
2 |  ('i', 1267),  
3 |  ('and', 1245),  
4 |  ('of', 1155),  
5 |  ('a', 816),  
6 |  ('to', 695),  
7 |  ('was', 552),  
8 |  ('in', 541),  
9 |  ('that', 443),  
10 | ('my', 440)]
```

```
1 | freqs = [freq for token, freq in vocab.token_freqs]  
2 | plt.plot(freqs)  
3 | plt.xlabel('token: x')  
4 | plt.ylabel('frequency: n(x)')  
5 | plt.xscale('log')  
6 | plt.yscale('log')
```



- We are on to something quite fundamental here: the word frequency decays rapidly in a well-defined way. After dealing with the first few words as exceptions, all the remaining words roughly follow a straight line on a log-log plot. This means that words satisfy Zipf's law, which states that the frequency n_i of the i^{th} most frequent word is:

$$n_i \propto \frac{1}{i^\alpha} \quad (13)$$
$$\Rightarrow \log n_i = -\alpha \log i + c,$$

- Where α is the exponent that characterizes the distribution and c is a constant.
- But what about the other word combinations, such as bigrams, trigrams, and beyond?

```
1 | bigram_tokens = [pair for pair in zip(corpus[:-1], corpus[1:])]  
2 | bigram_vocab = vocab(bigram_tokens)  
3 | bigram_vocab.token_freqs[:10]
```

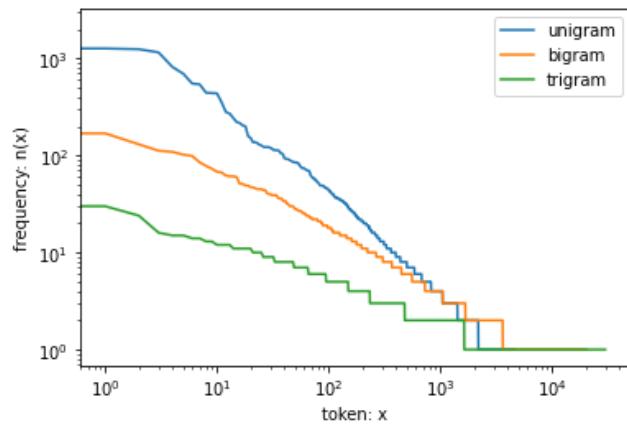
```
1  [(4, 1), 309),
2   (8, 1), 169),
3   (2, 12), 130),
4   (2, 7), 112),
5   (3, 1), 109),
6   (1, 19), 102),
7   (11, 7), 99),
8   (6, 1), 85),
9   (14, 2), 78),
10  (4, 5), 73)]
```

```
1  trigram_tokens = [triple for triple in zip(
2      corpus[:-2], corpus[1:-1], corpus[2:])]
3  trigram_vocab = Vocab(trigram_tokens)
4  trigram_vocab.token_freqs[:10]
```

```
1  [(1, 19, 71), 59),
2   (1, 19, 50), 30),
3   (1, 177, 65), 24),
4   (11, 63, 6), 16),
5   (11, 7, 5), 15),
6   (103, 3, 26), 15),
7   (63, 6, 13), 14),
8   (2, 107, 33), 14),
9   (2, 47, 1), 13),
10  (2, 110, 6), 13)]
```

```
1  bigram_freqs = [freq for token, freq in bigram_vocab.token_freqs]
2  trigram_freqs = [freq for token, freq in trigram_vocab.token_freqs]
3
4  plt.plot(freqs)
5  plt.plot(bigram_freqs)
6  plt.plot(trigram_freqs)
7  plt.xlabel('token: x')
8  plt.ylabel('frequency: n(x)')
9  plt.xscale('log')
10 plt.yscale('log')
11 plt.legend(['unigram', 'bigram', 'trigram'])
```

```
1  <matplotlib.legend.Legend at 0x24746d89ca0>
```



- This figure is quite exciting for a number of reasons. First, beyond unigram words, sequences of words also appear to be following Zipf's law, albeit with a smaller exponent α , depending on the sequence length.
- Second, the number of distinct n -grams is not that large. This gives us hope that there is quite a lot of structure in language. Third, many n -grams occur very rarely, which makes Laplace smoothing rather unsuitable for language modeling. Instead, we will use deep learning based models.

Long Sequence Data

- When sequences get too long to be processed by models all at once, we may wish to split such sequences for reading. Now let us describe general strategies. Before introducing the model, let us assume that we will use a neural network to train a language model, where the network processes a minibatch of sequences with predefined length, say n time steps, at a time. Now the question is how to read minibatches of features and labels at random.
- The figure below shows all the different ways to obtain subsequences from an original text sequence, where a token at each time step corresponds to a character.

the time machine by h g wells
the time machine by h g wells
the time machine by h g wells
the time machine by h g wells
the time machine by h g wells
the time machine by h g wells

Random Sampling

- However, if we pick just one offset, there is limited coverage of all the possible subsequences for training our network. Therefore, we can start with a random offset to partition a sequence to get both coverage and randomness.
- In random sampling, each example is a subsequence arbitrarily captured on the original long sequence. The subsequences from two adjacent random minibatches during iteration are not necessarily adjacent on the original sequence. For language modeling, the target is to predict the next token based on what tokens we have seen so far, hence the labels are the original sequence, shifted by one token.

```
1 def seq_data_iter_random(corpus, batch_size, num_steps): #@save
2     """Generate a minibatch of subsequences using random sampling."""
3     # Start with a random offset (inclusive of `num_steps - 1`) to partition a
4     # sequence
5     corpus = corpus[random.randint(0, num_steps - 1):]
6     # Subtract 1 since we need to account for labels
7     num_subseqs = (len(corpus) - 1) // num_steps
8     # The starting indices for subsequences of length `num_steps`
9     initial_indices = list(range(0, num_subseqs * num_steps, num_steps))
```

```

10 # In random sampling, the subsequences from two adjacent random
11 # minibatches during iteration are not necessarily adjacent on the
12 # original sequence
13 random.shuffle(initial_indices)
14
15 def data(pos):
16     # Return a sequence of length `num_steps` starting from `pos`
17     return corpus[pos: pos + num_steps]
18
19 num_batches = num_subseqs // batch_size
20 for i in range(0, batch_size * num_batches, batch_size):
21     # Here, `initial_indices` contains randomized starting indices for
22     # subsequences
23     initial_indices_per_batch = initial_indices[i: i + batch_size]
24     x = [data(j) for j in initial_indices_per_batch]
25     y = [data(j + 1) for j in initial_indices_per_batch]
26     yield torch.tensor(x), torch.tensor(y)

```

```

1 my_seq = list(range(35))
2 for X, Y in seq_data_iter_random(my_seq, batch_size=2, num_steps=5):
3     print('X: ', x, '\nY:', y)

```

```

1 X: tensor([[5, 6, 7, 8, 9],
2           [0, 1, 2, 3, 4]])
3 Y: tensor([[ 6,  7,  8,  9, 10],
4           [ 1,  2,  3,  4,  5]])
5 X: tensor([[20, 21, 22, 23, 24],
6           [15, 16, 17, 18, 19]])
7 Y: tensor([[21, 22, 23, 24, 25],
8           [16, 17, 18, 19, 20]])
9 X: tensor([[25, 26, 27, 28, 29],
10            [10, 11, 12, 13, 14]])
11 Y: tensor([[26, 27, 28, 29, 30],
12            [11, 12, 13, 14, 15]])

```

```

1 def seq_data_iter_sequential(corpus, batch_size, num_steps): #@save
2     """Generate a minibatch of subsequences using sequential partitioning."""
3     # Start with a random offset to partition a sequence
4     offset = random.randint(0, num_steps)
5     num_tokens = ((len(corpus) - offset - 1) // batch_size) * batch_size
6     xs = torch.tensor(corpus[offset: offset + num_tokens])
7     ys = torch.tensor(corpus[offset + 1: offset + 1 + num_tokens])
8     xs, ys = xs.reshape(batch_size, -1), ys.reshape(batch_size, -1)
9     num_batches = xs.shape[1] // num_steps
10    for i in range(0, num_steps * num_batches, num_steps):
11        x = xs[:, i: i + num_steps]
12        y = ys[:, i: i + num_steps]
13        yield x, y

```

```

1 for X, Y in seq_data_iter_sequential(my_seq, batch_size=2, num_steps=5):
2     print('X: ', x, '\nY:', y)

```

```

1 X: tensor([[ 1,  2,  3,  4,  5],
2           [17, 18, 19, 20, 21]])
3 Y: tensor([[ 2,  3,  4,  5,  6],
4           [18, 19, 20, 21, 22]])
5 X: tensor([[ 6,  7,  8,  9, 10],
6           [22, 23, 24, 25, 26]])
7 Y: tensor([[ 7,  8,  9, 10, 11],
8           [23, 24, 25, 26, 27]])
9 X: tensor([[11, 12, 13, 14, 15],
10          [27, 28, 29, 30, 31]])
11 Y: tensor([[12, 13, 14, 15, 16],
12          [28, 29, 30, 31, 32]])

```

```

1 class SeqDataLoader: #@save
2     """An iterator to load sequence data."""
3     def __init__(self, batch_size, num_steps, use_random_iter, max_tokens):
4         if use_random_iter:
5             self.data_iter_fn = seq_data_iter_random
6         else:
7             self.data_iter_fn = seq_data_iter_sequential
8         self.corpus, self.vocab = load_corpus_time_machine(max_tokens)
9         self.batch_size, self.num_steps = batch_size, num_steps
10
11     def __iter__(self):
12         return self.data_iter_fn(self.corpus, self.batch_size, self.num_steps)

```

```

1 def load_data_time_machine(batch_size, num_steps, #@save
2                             use_random_iter=False, max_tokens=10000):
3     """Return the iterator and the vocabulary of the time machine dataset."""
4     data_iter = SeqDataLoader(batch_size, num_steps, use_random_iter, max_tokens)
5     return data_iter, data_iter.vocab

```

Recurrent Neural Networks

- We have introduced n -gram models, where the conditional probability of word x_t at time step t only depends on the $n - 1$ previous words. If we want to incorporate the possible effect of words earlier than time step $t - (n - 1)$ on x_t , we need to increase n . However, the number of model parameters would also increase exponentially with it, as we need to store $|V|^n$ numbers for a vocabulary set V . Hence, rather than modeling $P(x_t | x_{t-1}, \dots, x_{t-n+1})$ it is preferable to use a latent variable model:

$$P(x_t | x_{t-1}, \dots, x_1) \approx P(x_t | h_{t-1}) \quad (14)$$

- where h_{t-1} is a hidden state (also known as a hidden variable) that stores the sequence information up to time step $t - 1$. In general, the hidden state at any time step t could be computed based on both the current input x_t and the previous hidden state h_{t-1} :

$$h_t = f(x_t, h_{t-1}) \quad (15)$$

- For a sufficiently powerful function f , the latent variable model is not an approximation. After all, h_t may simply store all the data it has observed so far. However, it could potentially make both computation and storage expensive.
- Recall that we have discussed hidden layers with hidden units. It is noteworthy that hidden layers and hidden states refer to two very different concepts. Hidden layers are, as explained, layers that are hidden from view on the path from input to output. Hidden states are technically speaking inputs to whatever we do at a given step, and they can only be computed by looking at data at previous time steps.
- Recurrent neural networks (RNNs) are neural networks with hidden states.
- Let us take a look at an MLP with a single hidden layer. Let the hidden layer's activation function be ϕ . Given a minibatch of examples $\mathbf{X} \in \mathbb{R}^{n \times d}$ with batch size n and d inputs, the hidden layer's output $\mathbf{H} \in \mathbb{R}^{n \times h}$ is calculated as:

$$\mathbf{H} = \phi(\mathbf{X}\mathbf{W}_{xh} + \mathbf{b}_h) \quad (16)$$

- where $\mathbf{W}_{xq} \in \mathbb{R}^{d \times h}$ is the weight parameter, and $\mathbf{b}_h \in \mathbb{R}^{1 \times h}$ is the bias parameter.
- Next, the hidden variable \mathbf{H} is used as the input of the output layer. The output layer is given by:

$$\mathbf{O} = \mathbf{H}\mathbf{W}_{hq} + \mathbf{b}_q \quad (17)$$

- where $\mathbf{O} \in \mathbb{R}^{n \times q}$ is the output variable, $\mathbf{W}_{hq} \in \mathbb{R}^{h \times q}$ is the weight parameter, and $\mathbf{b}_q \in \mathbb{R}^{1 \times q}$ is the bias parameter of the output layer. If it is a classification problem, we can use $\text{softmax}(\mathbf{O})$ to compute the probability distribution of the output categories.
- Assume that we have a minibatch of inputs $\mathbf{X}_t \in \mathbb{R}^{n \times d}$ at time step t . In other words, for a minibatch of n sequence examples, each row of \mathbf{X}_t corresponds to one example at time step t from the sequence. Next, denote by $\mathbf{H}_t \in \mathbb{R}^{n \times h}$ the hidden variable of time step t .
- Unlike the MLP, here we save the hidden variable \mathbf{H}_{t-1} from the previous time step and introduce a new weight parameter $\mathbf{W}_{hh} \in \mathbb{R}^{h \times h}$ to describe how to use the hidden variable of the previous time step in the current time step. Specifically, the calculation of the hidden variable of the current time step is determined by the input of the current time step together with the hidden variable of the previous time step:

$$\mathbf{H}_t = \phi(\mathbf{X}_t \mathbf{W}_{xh} + \mathbf{H}_{t-1} \mathbf{W}_{hh} + \mathbf{b}_h) \quad (18)$$

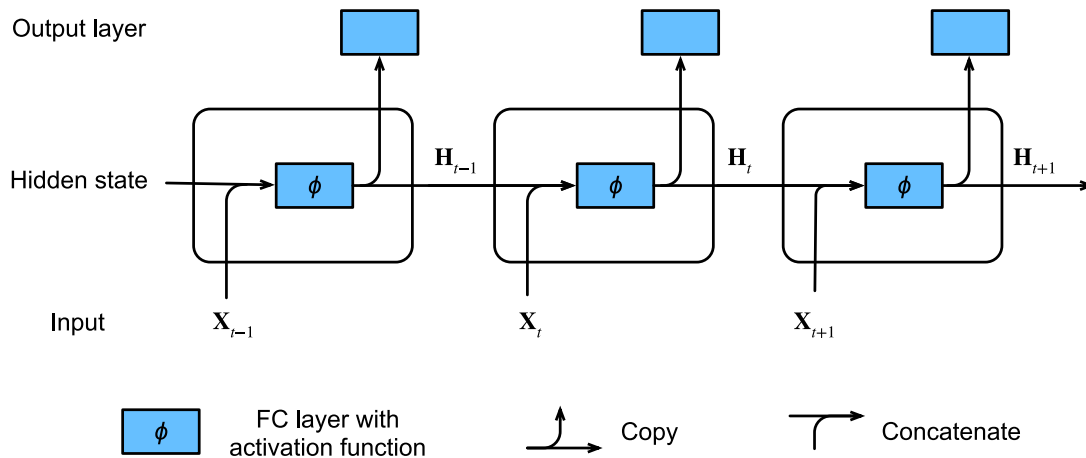
- Compared with:

$$\mathbf{H} = \phi(\mathbf{X} \mathbf{W}_{xh} + \mathbf{b}_h)$$

- this one adds one more term $\mathbf{H}_{t-1} \mathbf{W}_{hh}$ and thus instantiates $h_t = f(x_t, h_{t-1})$. From the relationship between hidden variables \mathbf{H}_t and \mathbf{H}_{t-1} of adjacent time steps, we know that these variables captured and retained the sequence's historical information up to their current time step, just like the state or memory of the neural network's current time step.
- Therefore, such a hidden variable is called a hidden state. Since the hidden state uses the same definition of the previous time step in the current time step, the computation is recurrent.
- Hence, neural networks with hidden states based on recurrent computation are named recurrent neural networks. Layers that perform the computation in RNNs are called recurrent layers.
- For time step, the output of the output layer of RNN is similar to the computation in the MLP:

$$\mathbf{O} = \mathbf{H} \mathbf{W}_{hq} + \mathbf{b}_q \quad (19)$$

- Parameters of the RNN include the weights $\mathbf{W}_{xh} \in \mathbb{R}^{d \times h}$, $\mathbf{W}_{hh} \in \mathbb{R}^{h \times h}$, and the bias $\mathbf{b}_h \in \mathbb{R}^{1 \times h}$ of the hidden layer, together with the weights $\mathbf{W}_{hq} \in \mathbb{R}^{h \times q}$ and the bias $\mathbf{b}_q \in \mathbb{R}^{1 \times q}$ of the output layer.



```

1 # A trick
2 X, w_xh = torch.normal(0, 1, (3, 1)), torch.normal(0, 1, (1, 4))
3 H, w_hh = torch.normal(0, 1, (3, 4)), torch.normal(0, 1, (4, 4))
4 print(torch.matmul(X, w_xh) + torch.matmul(H, w_hh))
5 torch.matmul(torch.cat((X, H), 1), torch.cat((w_xh, w_hh), 0))

```

```

1 tensor([[ 1.3129, -1.2245,  2.8187, -2.7141],
2         [-2.1283,  1.8983, -1.7776, -0.9191],
3         [-0.7369, -1.8983,  4.0983,  4.2379]])

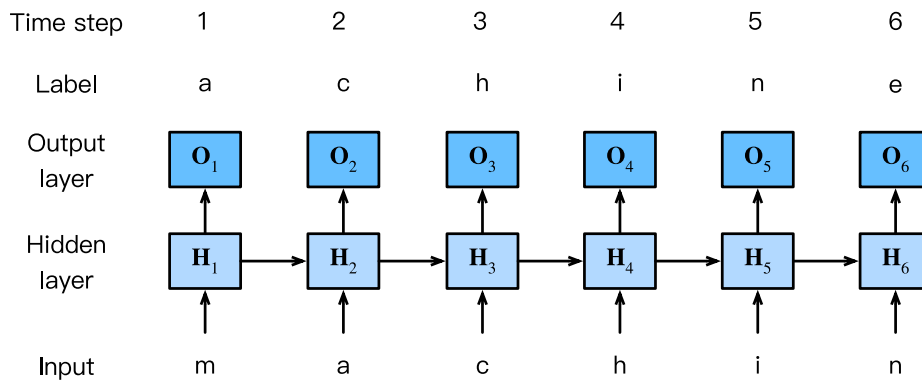
```

```

1 tensor([[ 1.3129, -1.2245,  2.8187, -2.7141],
2         [-2.1283,  1.8983, -1.7776, -0.9191],
3         [-0.7369, -1.8983,  4.0983,  4.2379]])

```

- Recall that for language modeling, we aim to predict the next token based on the current and past tokens, thus we shift the original sequence by one token as the labels. Bengio et al. first proposed to use a neural network for language modeling
- Let the minibatch size be one, and the sequence of the text be “machine”. To simplify training in subsequent sections, we tokenize text into characters rather than words and consider a character-level language model. The figure below demonstrates how to predict the next character based on the current and previous characters via an RNN for character-level language modeling.



A Neural Probabilistic Language Model

Yoshua Bengio
Réjean Ducharme
Pascal Vincent
Christian Jauvin

Département d'Informatique et Recherche Opérationnelle
Centre de Recherche Mathématiques
Université de Montréal, Montréal, Québec, Canada

BENGIOY@IRO.UMONTREAL.CA
 DUCHARME@IRO.UMONTREAL.CA
 VINCENTP@IRO.UMONTREAL.CA
 JAUVINC@IRO.UMONTREAL.CA

Editors: Jaz Kandola, Thomas Hofmann, Tomaso Poggio and John Shawe-Taylor

Abstract

A goal of statistical language modeling is to learn the joint probability function of sequences of words in a language. This is intrinsically difficult because of the **curse of dimensionality**: a word sequence on which the model will be tested is likely to be different from all the word sequences seen during training. Traditional but very successful approaches based on n-grams obtain generalization by concatenating very short overlapping sequences seen in the training set. We propose to fight the curse of dimensionality by **learning a distributed representation for words** which allows each training sentence to inform the model about an exponential number of semantically neighboring sentences. The model learns simultaneously (1) a distributed representation for each word along with (2) the probability function for word sequences, expressed in terms of these representations. Generalization is obtained because a sequence of words that has never been seen before gets high probability if it is made of words that are similar (in the sense of having a nearby representation) to words forming an already seen sentence. Training such large models (with millions of parameters) within a reasonable time is itself a significant challenge. We report on experiments using neural networks for the probability function, showing on two text corpora that the proposed approach significantly improves on state-of-the-art n-gram models, and that the proposed approach allows to take advantage of longer contexts.

Keywords: Statistical language modeling, artificial neural networks, distributed representation, curse of dimensionality

- Last, let us discuss about how to measure the language model quality, which will be used to evaluate our RNN-based models in the subsequent sections. One way is to check how surprising the text is. A good language model is able to predict with high-accuracy tokens that what we will see next. Consider the following continuations of the phrase “It is raining”, as proposed by different language models:

- "It is raining outside"
- "It is raining banana tree"
- "It is raining piouw;kcj pwepoiut"
- We might measure the quality of the model by computing the likelihood of the sequence. Unfortunately this is a number that is hard to understand and difficult to compare. After all, shorter sequences are much more likely to occur than the longer ones, hence evaluating the model on Tolstoy's magnum opus War and Peace will inevitably produce a much smaller likelihood than, say, on Saint-Exupery's novella The Little Prince. What is missing is the equivalent of an average.
- Information theory comes handy here. We have defined entropy, surprisal, and cross-entropy when we introduced the softmax regression. If we want to compress text, we can ask about predicting the next token given the current set of tokens. A better language model should allow us to predict the next token more accurately. Thus, it should allow us to spend fewer bits in compressing the sequence. So we can measure it by the cross-entropy loss averaged over all the tokens of a sequence:

$$\frac{1}{n} \sum_{t=1}^n -\log P(x_t | x_{t-1}, \dots, x_1) \quad (20)$$

- where P is given by a language model and x_t is the actual token observed at time step t from the sequence. This makes the performance on documents of different lengths comparable. For historical reasons, scientists in natural language processing prefer to use a quantity called perplexity. In a nutshell, it is the exponential of the equation above:

$$\exp \left(-\frac{1}{n} \sum_{t=1}^n \log P(x_t | x_{t-1}, \dots, x_1) \right) \quad (21)$$

- Perplexity can be best understood as the harmonic mean of the number of real choices that we have when deciding which token to pick next. Let us look at a number of cases:
 - In the best case scenario, the model always perfectly estimates the probability of the label token as 1. In this case the perplexity of the model is 1.
 - In the worst case scenario, the model always predicts the probability of the label token as 0. In this situation, the perplexity is positive infinity.
 - At the baseline, the model predicts a uniform distribution over all the available tokens of the vocabulary. In this case, the perplexity equals the number of unique tokens of the vocabulary. In fact, if we were to store the sequence without any compression, this would be the best we could do to encode it. Hence, this provides a nontrivial upper bound that any useful model must beat.

```

1 batch_size, num_steps = 32, 35
2 train_iter, vocab = load_data_time_machine(batch_size, num_steps)
3 print(F.one_hot(torch.tensor([0, 2])))
4 print("len(vocab) ->", len(vocab))
5 x = torch.arange(10).reshape((2, 5))
6 print(F.one_hot(X.T, 28).shape)

```

```

1 tensor([[1, 0, 0],
2         [0, 0, 1]])
3 len(vocab) -> 28
4 torch.Size([5, 2, 28])

```

```

1 def get_params(vocab_size, num_hiddens, device):
2     num_inputs = num_outputs = vocab_size
3
4     def normal(shape):
5         return torch.randn(size=shape, device=device) * 0.01
6
7     # Hidden layer parameters
8     w_xh = normal((num_inputs, num_hiddens))
9     w_hh = normal((num_hiddens, num_hiddens))
10    b_h = torch.zeros(num_hiddens, device=device)
11    # Output layer parameters
12    w_hq = normal((num_hiddens, num_outputs))

```

```

13     b_q = torch.zeros(num_outputs, device=device)
14     # Attach gradients
15     params = [w_xh, w_hh, b_h, w_hq, b_q]
16     for param in params:
17         param.requires_grad_(True)
18     return params

```

Build RNN Model

```

1 def init_rnn_state(batch_size, num_hiddens, device):
2     return (torch.zeros((batch_size, num_hiddens), device=device), )
3
4 def rnn_fw(inputs, state, params):
5     # Here `inputs` shape: (`num_steps`, `batch_size`, `vocab_size`)
6     w_xh, w_hh, b_h, w_hq, b_q = params
7     H, = state
8     outputs = []
9     # Shape of `x`: (`batch_size`, `vocab_size`)
10    for x in inputs:
11        H = torch.tanh(torch.mm(x, w_xh) + torch.mm(H, w_hh) + b_h)
12        Y = torch.mm(H, w_hq) + b_q
13        outputs.append(Y)
14    return torch.cat(outputs, dim=0), (H,)

```

```

1 class RNN:
2     """A RNN Model implemented from scratch."""
3     def __init__(self, vocab_size, num_hiddens, device, get_params, init_state, forward_fn):
4         self.vocab_size, self.num_hiddens = vocab_size, num_hiddens
5         self.params = get_params(vocab_size, num_hiddens, device)
6         self.init_state, self.forward_fn = init_state, forward_fn
7
8     def __call__(self, x, state):
9         # print(x.T)
10        # print(x.shape)
11        x = F.one_hot(x.T, self.vocab_size).type(torch.float32)
12        return self.forward_fn(x, state, self.params)
13
14    def begin_state(self, batch_size, device):
15        return self.init_state(batch_size, self.num_hiddens, device)

```

```

1 num_hiddens = 512
2 device = "cuda" if torch.cuda.is_available() else "cpu"
3 print(torch.cuda.get_arch_list(), device)
4 rnn = RNN(len(vocab), num_hiddens, device, get_params, init_rnn_state, rnn_fw)
5 state = rnn.begin_state(X.shape[0], device)
6 print(isinstance(state, tuple), type(state))
7 Y, new_state = rnn(X.to(device), state)
8 print(Y.shape, len(new_state), new_state[0].shape)

```

```

1 ['sm_37', 'sm_50', 'sm_60', 'sm_61', 'sm_70', 'sm_75', 'sm_80', 'sm_86', 'compute_37'] cuda
2 True <class 'tuple'>
3 torch.Size([10, 28]) 1 torch.Size([2, 512])

```

```

1 def predict(prefix, num_preds, net, vocab, device):
2     """Generate new characters following the `prefix`."""
3     state = net.begin_state(batch_size=1, device=device)
4     outputs = [vocab[prefix[0]]]
5     get_input = lambda: torch.tensor([outputs[-1]], device=device).reshape((1, 1))
6     for y in prefix[1:]: # warm-up period
7         _, state = net(get_input(), state)
8         outputs.append(vocab[y])
9     for _ in range(num_preds): # Predict `num_preds` steps
10        y, state = net(get_input(), state)
11        outputs.append(int(y.argmax(dim=1).reshape(1)))
12    return ''.join([vocab.idx_to_token[i] for i in outputs])

```

```

1 predict('time traveller ', 10, rnn, vocab, device)

```

```

1 'time traveller oidxoqrjz '

```

Gradient Clipping

- For a sequence of length T , we compute the gradients over these T time steps in an iteration, which results in a chain of matrix-products with length $O(T)$ during backpropagation. It might result in numerical instability, e.g., the gradients may either explode or vanish, when T is large. Therefore, RNN models often need extra help to stabilize the training.
- Generally speaking, when solving an optimization problem, we take update steps for the model parameter, say in the vector form \mathbf{x} , in the direction of the negative gradient \mathbf{g} on a minibatch. For example, with $\eta > 0$ as the learning rate, in one iteration we update \mathbf{x} as $\mathbf{x} - \eta \mathbf{g}$. Let us further assume that the objective function f is well behaved, say, Lipschitz continuous with constant L . That is to say, for any \mathbf{a} and \mathbf{b} we have:

$$|f(\mathbf{a}) - f(\mathbf{b})| \leq L \|\mathbf{a} - \mathbf{b}\| \quad (22)$$

- In this case we can safely assume that if we update the parameter vector by $\eta \mathbf{g}$, then:

$$|f(\mathbf{x}) - f(\mathbf{x} - \eta \mathbf{g})| \leq L \eta \|\mathbf{g}\| \quad (23)$$

- which means that we will not observe a change by more than $\eta \|\mathbf{g}\|$. This is both a curse and a blessing. On the curse side, it limits the speed of making progress; whereas on the blessing side, it limits the extent to which things can go wrong if we move in the wrong direction.
- Sometimes the gradients can be quite large and the optimization algorithm may fail to converge. We could address this by reducing the learning rate η . But what if we only rarely get large gradients? In this case such an approach may appear entirely unwarranted. One popular alternative is to clip the gradient \mathbf{g} by projecting them back to a ball of a given radius, say θ via:

$$g \leftarrow \min \left(1, \frac{\theta}{\|\mathbf{g}\|} \right) \mathbf{g} \quad (24)$$

```

1 def grad_clipping(net, theta): #@save
2     """Clip the gradient."""
3     if isinstance(net, nn.Module):
4         params = [p for p in net.parameters() if p.requires_grad]
5     else:
6         params = net.params
7     norm = torch.sqrt(sum(torch.sum((p.grad ** 2)) for p in params))
8     if norm > theta:
9         for param in params:
10            param.grad[:] *= theta / norm

```

```

1 import time
2 def train_epoch(net, train_iter, loss, optimizer, device, use_random_iter):
3     """Train a net within one epoch"""
4     state, start_time = None, time.perf_counter()
5     metric = [0, 0]
6     loss_sum = 0

```

```

7   size = 0
8   for X, y in train_iter:
9       size += 1
10  for X, y in train_iter:
11      if state is None or use_random_iter:
12          # Initialize `state` when either it is the first iteration or
13          # using random sampling
14          state = net.begin_state(batch_size=X.shape[0], device=device)
15      else:
16          if isinstance(net, nn.Module) and not isinstance(state, tuple):
17              # `state` is a tensor for `nn.GRU`
18              state.detach_()
19          else:
20              # `state` is a tuple of tensors for `nn.LSTM` and
21              # for our custom scratch implementation
22              for s in state:
23                  s.detach_()
24      y = Y.T.reshape(-1)
25      X, y = X.to(device), y.to(device)
26      y_hat, state = net(X, state)
27      l = loss(y_hat, y.long()).mean()
28      loss_sum += l
29      if isinstance(optimizer, torch.optim.Optimizer):
30          optimizer.zero_grad()
31          l.backward()
32          grad_clipping(net, 1)
33          optimizer.step()
34      else:
35          l.backward()
36          grad_clipping(net, 1)
37          optimizer(batch_size=1) # Since the `mean` function has been invoked
38      metric[0] += l * y.numel()
39      metric[1] += y.numel()
40  return math.exp(metric[0]/metric[1]), metric[1]/(time.perf_counter() - start_time),
    loss_sum/size

```

```

1  def sgd(params, lr, batch_size):
2      """Minibatch stochastic gradient descent. Defined in :numref:`sec_linear_scratch`"""
3      with torch.no_grad():
4          for param in params:
5              param -= lr * param.grad / batch_size
6              param.grad.zero_()

```

```

1  def train_rnn(net, train_iter, vocab, lr, num_epochs, device, use_random_iter=False):
2      """Train a model"""
3      loss = nn.CrossEntropyLoss()
4      loss_list = []
5      # Initialize
6      if isinstance(net, nn.Module):
7          optimizer = torch.optim.SGD(net.parameters(), lr)
8      else:
9          optimizer = lambda batch_size: sgd(net.params, lr, batch_size)
10     # Train and predict
11     for epoch in range(num_epochs):
12         ppl, speed, l = train_epoch(net, train_iter, loss, optimizer, device, use_random_iter)
13         loss_list.append(l.cpu().detach().numpy().sum())
14         if (epoch + 1) % 25 == 0 or epoch == 0:
15             print(predict('time traveller', 50, net, vocab, device))
16             print("Epoch ->", epoch + 1, "\t\t\tLoss ->", l)
17         print('perplexity ->', ppl, "-*-*-*-", speed, 'tokens/sec on', device)
18         print(predict('time traveller', 50, net, vocab, device))
19         print(predict('traveller', 50 + 5, net, vocab, device))
20     return loss_list

```

```

1 def plot_losses(loss_list):
2     plt.figure(figsize=(20, 10))
3     plt.plot(loss_list, '-bx')
4     plt.xlabel('epoch')
5     plt.ylabel('loss')
6     plt.title('Loss vs. No. of epochs')
7     plt.show()

```

```

1 num_epochs, lr = 2000, 0.5
2 loss_list = train_rnn(rnn, train_iter, vocab, lr, num_epochs, device)

```

```

1 time traveller
2 Epoch -> 1          Loss -> tensor(3.2657, device='cuda:0', grad_fn=<DivBackward0>)
3 time traveller the the the the the the the the the the the t
4 Epoch -> 25         Loss -> tensor(2.4908, device='cuda:0', grad_fn=<DivBackward0>)
5 time traveller the the the the the the the the the the the
6 Epoch -> 50         Loss -> tensor(2.2883, device='cuda:0', grad_fn=<DivBackward0>)
7 time traveller the the the the the the the the the the the
8 Epoch -> 75         Loss -> tensor(2.1909, device='cuda:0', grad_fn=<DivBackward0>)
9 time traveller and the the the the the the the the the the t
10 Epoch -> 100        Loss -> tensor(2.1230, device='cuda:0', grad_fn=<DivBackward0>)
11 time traveller and and and and and and and and and and and
12 Epoch -> 125        Loss -> tensor(2.0727, device='cuda:0', grad_fn=<DivBackward0>)
13 time traveller and the the the the the the the the the the t
14 Epoch -> 150        Loss -> tensor(2.0153, device='cuda:0', grad_fn=<DivBackward0>)
15 time traveller and and the the the the the the the the the t
16 Epoch -> 175        Loss -> tensor(1.9479, device='cuda:0', grad_fn=<DivBackward0>)
17 time traveller at in the ghat and have and the that these the th
18 Epoch -> 200        Loss -> tensor(1.8707, device='cuda:0', grad_fn=<DivBackward0>)
19 time traveller ourd and the time traveller ourd and the time tra
20 Epoch -> 225        Loss -> tensor(1.7279, device='cuda:0', grad_fn=<DivBackward0>)
21 time traveller of the thing s arall or the fill tracell same the
22 Epoch -> 250        Loss -> tensor(1.4753, device='cuda:0', grad_fn=<DivBackward0>)
23 time traveller bet and the time traveller dewe mith his in the g
24 Epoch -> 275        Loss -> tensor(1.2242, device='cuda:0', grad_fn=<DivBackward0>)
25 time traveller pereat in tha file wather sacco the grave thavelle
26 Epoch -> 300        Loss -> tensor(0.9787, device='cuda:0', grad_fn=<DivBackward0>)
27 time traveller ffre dine ti noous it in ass if eredinger as wa h
28 Epoch -> 325        Loss -> tensor(0.7369, device='cuda:0', grad_fn=<DivBackward0>)
29 time traveller come the bouthe buine ffreew yhr gereatthred hit
30 Epoch -> 350        Loss -> tensor(0.5201, device='cuda:0', grad_fn=<DivBackward0>)
31 time traveller came back anduiney filbes sthey ubuin expaysedce
32 Epoch -> 375        Loss -> tensor(0.3837, device='cuda:0', grad_fn=<DivBackward0>)
33 time traveller come the gat ee that ale ingtome wathe one ur ato
34 Epoch -> 400        Loss -> tensor(0.2849, device='cuda:0', grad_fn=<DivBackward0>)
35 time traveller cald twhe bogit bote mod t is the three dimension
36 Epoch -> 425        Loss -> tensor(0.2002, device='cuda:0', grad_fn=<DivBackward0>)
37 time traveller after the pauserequired for the proper assimilati
38 Epoch -> 450        Loss -> tensor(0.1134, device='cuda:0', grad_fn=<DivBackward0>)
39 time travellerif s aga beree by a couve buthed heare hea sime tr
40 Epoch -> 475        Loss -> tensor(0.0767, device='cuda:0', grad_fn=<DivBackward0>)
41 time traveller for so it will be convenient to speak of himwas e
42 Epoch -> 500        Loss -> tensor(0.0957, device='cuda:0', grad_fn=<DivBackward0>)
43 time traveller you can show black is white by argument said filby
44 Epoch -> 525        Loss -> tensor(0.0795, device='cuda:0', grad_fn=<DivBackward0>)
45 time traveller you can show black is white by argument said filby
46 Epoch -> 550        Loss -> tensor(0.0562, device='cuda:0', grad_fn=<DivBackward0>)
47 time traveller you can show black is white by argument said filby
48 Epoch -> 575        Loss -> tensor(0.0747, device='cuda:0', grad_fn=<DivBackward0>)
49 time traveller for so it will be convenient to speak of himwas e
50 Epoch -> 600        Loss -> tensor(0.0515, device='cuda:0', grad_fn=<DivBackward0>)
51 time traveller for so it will be convenient to speak of himwas e
52 Epoch -> 625        Loss -> tensor(0.0732, device='cuda:0', grad_fn=<DivBackward0>)
53 time traveller you can show black is white by argument said filby
54 Epoch -> 650        Loss -> tensor(0.0447, device='cuda:0', grad_fn=<DivBackward0>)
55 time traveller you can show black is white by argument said filby

```



```
56 Epoch -> 675          Loss -> tensor(0.0458, device='cuda:0', grad_fn=<DivBackward0>)
57 time traveller you can show black is white by argument said filby
58 Epoch -> 700          Loss -> tensor(0.0584, device='cuda:0', grad_fn=<DivBackward0>)
59 time traveller with a slight accession of cheerfulness really thi
60 Epoch -> 725          Loss -> tensor(0.0574, device='cuda:0', grad_fn=<DivBackward0>)
61 time traveller you can show black is white by argument said filby
62 Epoch -> 750          Loss -> tensor(0.0296, device='cuda:0', grad_fn=<DivBackward0>)
63 time traveller you can show black is white by argument said filby
64 Epoch -> 775          Loss -> tensor(0.0593, device='cuda:0', grad_fn=<DivBackward0>)
65 time traveller you can show black is white by argument said filby
66 Epoch -> 800          Loss -> tensor(0.0437, device='cuda:0', grad_fn=<DivBackward0>)
67 time traveller for so it will be convenient to speak of him was e
68 Epoch -> 825          Loss -> tensor(0.0480, device='cuda:0', grad_fn=<DivBackward0>)
69 time traveller for so it will be convenient to speak of him was e
70 Epoch -> 850          Loss -> tensor(0.0452, device='cuda:0', grad_fn=<DivBackward0>)
71 time traveller for so it will be convenient to speak of him was e
72 Epoch -> 875          Loss -> tensor(0.0394, device='cuda:0', grad_fn=<DivBackward0>)
73 time traveller you can show black is white by argument said filby
74 Epoch -> 900          Loss -> tensor(0.0419, device='cuda:0', grad_fn=<DivBackward0>)
75 time traveller for so it will be convenient to speak of him was e
76 Epoch -> 925          Loss -> tensor(0.0409, device='cuda:0', grad_fn=<DivBackward0>)
77 time traveller you can show black is white by argument said filby
78 Epoch -> 950          Loss -> tensor(0.0304, device='cuda:0', grad_fn=<DivBackward0>)
79 time traveller you can show black is white by argument said filby
80 Epoch -> 975          Loss -> tensor(0.0319, device='cuda:0', grad_fn=<DivBackward0>)
81 time traveller with a slight accession of cheerfulness really thi
82 Epoch -> 1000         Loss -> tensor(0.0342, device='cuda:0', grad_fn=<DivBackward0>)
83 time traveller with a slight accession of cheerfulness really thi
84 Epoch -> 1025         Loss -> tensor(0.0316, device='cuda:0', grad_fn=<DivBackward0>)
85 time traveller with a slight accession of cheerfulness really thi
86 Epoch -> 1050         Loss -> tensor(0.0253, device='cuda:0', grad_fn=<DivBackward0>)
87 time traveller with a slight accession of cheerfulness really thi
88 Epoch -> 1075         Loss -> tensor(0.0296, device='cuda:0', grad_fn=<DivBackward0>)
89 time traveller for so it will be convenient to speak of him was e
90 Epoch -> 1100         Loss -> tensor(0.0252, device='cuda:0', grad_fn=<DivBackward0>)
91 time traveller for so it will be convenient to speak of him was e
92 Epoch -> 1125         Loss -> tensor(0.0311, device='cuda:0', grad_fn=<DivBackward0>)
93 time traveller you can show black is white by argument said filby
94 Epoch -> 1150         Loss -> tensor(0.0289, device='cuda:0', grad_fn=<DivBackward0>)
95 time traveller you can show black is white by argument said filby
96 Epoch -> 1175         Loss -> tensor(0.0285, device='cuda:0', grad_fn=<DivBackward0>)
97 time traveller you can show black is white by argument said filby
98 Epoch -> 1200         Loss -> tensor(0.0186, device='cuda:0', grad_fn=<DivBackward0>)
99 time traveller with a slight accession of cheerfulness really thi
10 Epoch -> 1225         Loss -> tensor(0.0257, device='cuda:0', grad_fn=<DivBackward0>)
10 time traveller you can show black is white by argument said filby
10 Epoch -> 1250         Loss -> tensor(0.0237, device='cuda:0', grad_fn=<DivBackward0>)
10 time traveller you can show black is white by argument said filby
10 Epoch -> 1275         Loss -> tensor(0.0222, device='cuda:0', grad_fn=<DivBackward0>)
10 time traveller you can show black is white by argument said filby
10 Epoch -> 1300         Loss -> tensor(0.0246, device='cuda:0', grad_fn=<DivBackward0>)
10 time traveller tered wit sard ily thaverisy cuvelyer io is aterly
10 Epoch -> 1325         Loss -> tensor(0.1075, device='cuda:0', grad_fn=<DivBackward0>)
10 time traveller you can show black is white by argument said filby
10 Epoch -> 1350         Loss -> tensor(0.0585, device='cuda:0', grad_fn=<DivBackward0>)
10 time traveller for so it will be convenient to speak of him was e
11 Epoch -> 1375         Loss -> tensor(0.0340, device='cuda:0', grad_fn=<DivBackward0>)
12 time traveller you can show black is white by argument said filby
13 Epoch -> 1400         Loss -> tensor(0.0300, device='cuda:0', grad_fn=<DivBackward0>)
14 time traveller for so it will be convenient to speak of him was e
15 Epoch -> 1425         Loss -> tensor(0.0257, device='cuda:0', grad_fn=<DivBackward0>)
16 time traveller you can show black is white by argument said filby
17 Epoch -> 1450         Loss -> tensor(0.0223, device='cuda:0', grad_fn=<DivBackward0>)
18 time traveller for so it will be convenient to speak of him was e
19 Epoch -> 1475         Loss -> tensor(0.0307, device='cuda:0', grad_fn=<DivBackward0>)
10 time traveller you can show black is white by argument said filby
12 Epoch -> 1500         Loss -> tensor(0.0224, device='cuda:0', grad_fn=<DivBackward0>)
12 time traveller you can show black is white by argument said filby
```

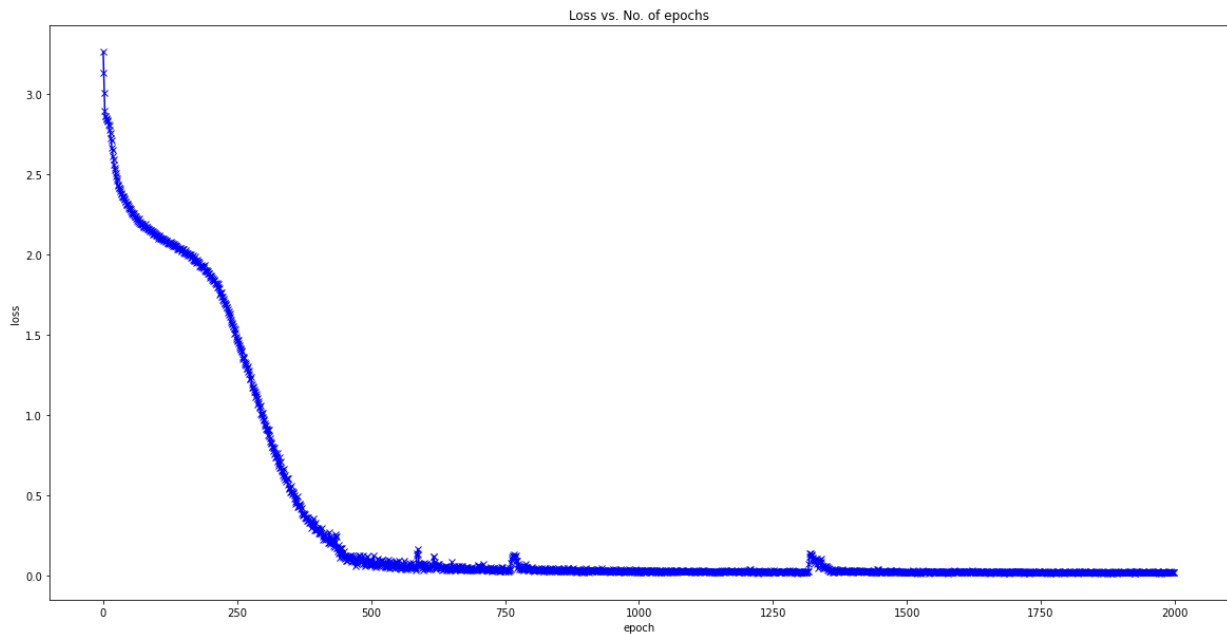


```

12 Epoch -> 1525          Loss -> tensor(0.0186, device='cuda:0', grad_fn=<DivBackward0>)
12 time traveller with a slight accession of cheerfulness really thi
13 Epoch -> 1550          Loss -> tensor(0.0231, device='cuda:0', grad_fn=<DivBackward0>)
13 time traveller for so it will be convenient to speak of himwas e
12 Epoch -> 1575          Loss -> tensor(0.0302, device='cuda:0', grad_fn=<DivBackward0>)
13 time traveller with a slight accession of cheerfulness really thi
19 Epoch -> 1600          Loss -> tensor(0.0251, device='cuda:0', grad_fn=<DivBackward0>)
10 time traveller you can show black is white by argument said filby
13 Epoch -> 1625          Loss -> tensor(0.0239, device='cuda:0', grad_fn=<DivBackward0>)
12 time traveller you can show black is white by argument said filby
13 Epoch -> 1650          Loss -> tensor(0.0247, device='cuda:0', grad_fn=<DivBackward0>)
13 time traveller for so it will be convenient to speak of himwas e
15 Epoch -> 1675          Loss -> tensor(0.0246, device='cuda:0', grad_fn=<DivBackward0>)
16 time traveller for so it will be convenient to speak of himwas e
13 Epoch -> 1700          Loss -> tensor(0.0238, device='cuda:0', grad_fn=<DivBackward0>)
13 time traveller you can show black is white by argument said filby
19 Epoch -> 1725          Loss -> tensor(0.0219, device='cuda:0', grad_fn=<DivBackward0>)
10 time traveller for so it will be convenient to speak of himwas e
14 Epoch -> 1750          Loss -> tensor(0.0243, device='cuda:0', grad_fn=<DivBackward0>)
12 time traveller for so it will be convenient to speak of himwas e
13 Epoch -> 1775          Loss -> tensor(0.0247, device='cuda:0', grad_fn=<DivBackward0>)
14 time traveller you can show black is white by argument said filby
15 Epoch -> 1800          Loss -> tensor(0.0193, device='cuda:0', grad_fn=<DivBackward0>)
16 time traveller for so it will be convenient to speak of himwas e
17 Epoch -> 1825          Loss -> tensor(0.0228, device='cuda:0', grad_fn=<DivBackward0>)
13 time traveller you can show black is white by argument said filby
19 Epoch -> 1850          Loss -> tensor(0.0209, device='cuda:0', grad_fn=<DivBackward0>)
10 time traveller you can show black is white by argument said filby
15 Epoch -> 1875          Loss -> tensor(0.0162, device='cuda:0', grad_fn=<DivBackward0>)
13 time traveller for so it will be convenient to speak of himwas e
13 Epoch -> 1900          Loss -> tensor(0.0210, device='cuda:0', grad_fn=<DivBackward0>)
14 time traveller with a slight accession of cheerfulness really thi
15 Epoch -> 1925          Loss -> tensor(0.0200, device='cuda:0', grad_fn=<DivBackward0>)
16 time traveller for so it will be convenient to speak of himwas e
13 Epoch -> 1950          Loss -> tensor(0.0226, device='cuda:0', grad_fn=<DivBackward0>)
13 time traveller for so it will be convenient to speak of himwas e
19 Epoch -> 1975          Loss -> tensor(0.0246, device='cuda:0', grad_fn=<DivBackward0>)
10 time traveller for so it will be convenient to speak of himwas e
16 Epoch -> 2000          Loss -> tensor(0.0244, device='cuda:0', grad_fn=<DivBackward0>)
13 perplexity -> 1.0246825922988203 -*- 110833.61474677507 tokens/sec on cuda
13 time traveller for so it will be convenient to speak of himwas e
16 traveller with a slight accession of cheerfulness really this is

```

```
1 | plot_losses(loss_list)
```



Use torch.nn.RNN()

```

1 rnn_layer = nn.RNN(len(vocab), num_hiddens)
2 init_state = torch.zeros((1, batch_size, num_hiddens))
3 print(init_state.shape)
4 X = torch.rand(size=(num_steps, batch_size, len(vocab)))
5 Y, state_new = rnn_layer(X, init_state)
6 print(Y.shape)
7 print(state_new.shape)

```

```

1 torch.Size([1, 32, 512])
2 torch.Size([35, 32, 512])
3 torch.Size([1, 32, 512])

```

```

1 class RNNModel(nn.Module):
2     """The RNN model."""
3     def __init__(self, rnn_layer, vocab_size, **kwargs):
4         super(RNNModel, self).__init__(**kwargs)
5         self.rnn = rnn_layer
6         self.vocab_size = vocab_size
7         self.num_hiddens = self.rnn.hidden_size
8         # If the RNN is bidirectional (to be introduced later), `num_directions` should be 2, else
9         # it should be 1.
10        if not self.rnn.bidirectional:
11            self.num_directions = 1
12            self.linear = nn.Linear(self.num_hiddens, self.vocab_size)
13        else:
14            self.num_directions = 2
15            self.linear = nn.Linear(self.num_hiddens * 2, self.vocab_size)
16
17        def forward(self, inputs, state):
18            X = F.one_hot(inputs.T.long(), self.vocab_size)
19            X = X.to(torch.float32)
20            Y, state = self.rnn(X, state)
21            # The fully connected layer will first change the shape of `Y` to (`num_steps` *
22            # `batch_size`, `num_hiddens`).
23            # Its output shape is (`num_steps` * `batch_size`, `vocab_size`).
24            output = self.linear(Y.reshape((-1, Y.shape[-1])))
25            return output, state

```

```

25     def begin_state(self, device, batch_size=1):
26         if not isinstance(self.rnn, nn.LSTM):
27             # `nn.GRU` takes a tensor as hidden state
28             return torch.zeros((self.num_directions * self.rnn.num_layers, batch_size,
self.num_hiddens), device=device)
29         else:
30             # `nn.LSTM` takes a tuple of hidden states
31             return (torch.zeros((self.num_directions * self.rnn.num_layers, batch_size,
self.num_hiddens), device=device),
32                     torch.zeros((self.num_directions * self.rnn.num_layers, batch_size,
self.num_hiddens), device=device))

```

Use torch.nn.RNN()

```

1 rnn_net = RNNModel(rnn_layer, vocab_size=len(vocab))
2 rnn_net = rnn_net.to(device)
3 predict('time traveller', 10, rnn_net, vocab, device)

```

```

1 'time travellerjkkeejeeeee'

```

```

1 num_epochs, lr = 2000, 0.5
2 loss_list = train_rnn(rnn_net, train_iter, vocab, lr, num_epochs, device)

```

```

1 time travellerreeeeeeeeeeeeeeeeeeeeeeeeeeeeeeeeeeeeeeeeeeeeeeeeeeeeee
2 Epoch -> 1          Loss -> tensor(3.0871, device='cuda:0', grad_fn=<DivBackward0>)
3 time traveller and and and and and and and and and and and
4 Epoch -> 25         Loss -> tensor(2.3187, device='cuda:0', grad_fn=<DivBackward0>)
5 time traveller the the the the the the the the the the the t
6 Epoch -> 50         Loss -> tensor(2.1041, device='cuda:0', grad_fn=<DivBackward0>)
7 time traveller and the the thas ans of the that in thave that ma
8 Epoch -> 75         Loss -> tensor(2.0062, device='cuda:0', grad_fn=<DivBackward0>)
9 time traveller mancenis thithed and har ghes anoul that ical s
10 Epoch -> 100        Loss -> tensor(1.9162, device='cuda:0', grad_fn=<DivBackward0>)
11 time traveller and ar anow youre ion ano stowe this loulhestire
12 Epoch -> 125        Loss -> tensor(1.6112, device='cuda:0', grad_fn=<DivBackward0>)
13 time traveller asting wio nesther if that ion is four whath le w
14 Epoch -> 150        Loss -> tensor(1.1728, device='cuda:0', grad_fn=<DivBackward0>)
15 time traveller ourde wor andthreat mave at carsome tha e and a m
16 Epoch -> 175        Loss -> tensor(0.7931, device='cuda:0', grad_fn=<DivBackward0>)
17 time traveller fred thes ather me sime traveller ags alw you but
18 Epoch -> 200        Loss -> tensor(0.5358, device='cuda:0', grad_fn=<DivBackward0>)
19 time traveller proceed d any t me time soa laves atwer and dome
20 Epoch -> 225        Loss -> tensor(0.3615, device='cuda:0', grad_fn=<DivBackward0>)
21 time traveller fur sowe ther tase phe lanse was axnot mask feorl
22 Epoch -> 250        Loss -> tensor(0.2797, device='cuda:0', grad_fn=<DivBackward0>)
23 time traveller for so it will be convenient to speak of himwas e
24 Epoch -> 275        Loss -> tensor(0.1188, device='cuda:0', grad_fn=<DivBackward0>)
25 time traveller fur no in and tie yor aclof absthrensther so deas
26 Epoch -> 300        Loss -> tensor(0.1235, device='cuda:0', grad_fn=<DivBackward0>)
27 time traveller fur ai in way h ao li g ou spas fhine bxe hat fof
28 Epoch -> 325        Loss -> tensor(0.1052, device='cuda:0', grad_fn=<DivBackward0>)
29 time traveller sure an hor aiman at move atoul ink in the glyere
30 Epoch -> 350        Loss -> tensor(0.0734, device='cuda:0', grad_fn=<DivBackward0>)
31 time traveller fit so in wall genely in the filbew of en as alo
32 Epoch -> 375        Loss -> tensor(0.0756, device='cuda:0', grad_fn=<DivBackward0>)
33 time traveller for so it will be convenient to speak of himwas e
34 Epoch -> 400        Loss -> tensor(0.0666, device='cuda:0', grad_fn=<DivBackward0>)
35 time traveller you can show black is white by argument said filby
36 Epoch -> 425        Loss -> tensor(0.0480, device='cuda:0', grad_fn=<DivBackward0>)
37 time traveller you can show black is white by argument said filby
38 Epoch -> 450        Loss -> tensor(0.0659, device='cuda:0', grad_fn=<DivBackward0>)

```

```
39 time traveller fith o inilitanig sury is of spe traveller you can
40 Epoch -> 475 Loss -> tensor(0.0572, device='cuda:0', grad_fn=<DivBackward0>)
41 time traveller for so it will be convenient to speak of him was e
42 Epoch -> 500 Loss -> tensor(0.0446, device='cuda:0', grad_fn=<DivBackward0>)
43 time traveller you can show black is white by argument said filby
44 Epoch -> 525 Loss -> tensor(0.0524, device='cuda:0', grad_fn=<DivBackward0>)
45 time traveller you can show black is white by argument said filby
46 Epoch -> 550 Loss -> tensor(0.0318, device='cuda:0', grad_fn=<DivBackward0>)
47 time traveller you can show black is white by argument said filby
48 Epoch -> 575 Loss -> tensor(0.0538, device='cuda:0', grad_fn=<DivBackward0>)
49 time traveller for so it will be convenient to speak of him was e
50 Epoch -> 600 Loss -> tensor(0.0499, device='cuda:0', grad_fn=<DivBackward0>)
51 time traveller for so it will be convenient to speak of him was e
52 Epoch -> 625 Loss -> tensor(0.0451, device='cuda:0', grad_fn=<DivBackward0>)
53 time traveller for so it will be convenient to speak of him was e
54 Epoch -> 650 Loss -> tensor(0.0375, device='cuda:0', grad_fn=<DivBackward0>)
55 time traveller you can show black is white by argument said filby
56 Epoch -> 675 Loss -> tensor(0.0415, device='cuda:0', grad_fn=<DivBackward0>)
57 time traveller you can show black is white by argument said filby
58 Epoch -> 700 Loss -> tensor(0.0381, device='cuda:0', grad_fn=<DivBackward0>)
59 time traveller for so it will be convenient to speak of him was e
60 Epoch -> 725 Loss -> tensor(0.0357, device='cuda:0', grad_fn=<DivBackward0>)
61 time traveller fith o inifatreyy ou tainve gab ta t otel ane ofm
62 Epoch -> 750 Loss -> tensor(0.0438, device='cuda:0', grad_fn=<DivBackward0>)
63 time traveller for so it will be convenient to speak of him was e
64 Epoch -> 775 Loss -> tensor(0.0577, device='cuda:0', grad_fn=<DivBackward0>)
65 time traveller you can show black is white by argument said filby
66 Epoch -> 800 Loss -> tensor(0.0488, device='cuda:0', grad_fn=<DivBackward0>)
67 time traveller for so it will be convenient to speak of him was e
68 Epoch -> 825 Loss -> tensor(0.0362, device='cuda:0', grad_fn=<DivBackward0>)
69 time traveller for so it will be convenient to speak of him was e
70 Epoch -> 850 Loss -> tensor(0.0292, device='cuda:0', grad_fn=<DivBackward0>)
71 time traveller with a slight accession of cheerfulness really thi
72 Epoch -> 875 Loss -> tensor(0.0297, device='cuda:0', grad_fn=<DivBackward0>)
73 time traveller for so it will be convenient to speak of him was e
74 Epoch -> 900 Loss -> tensor(0.0292, device='cuda:0', grad_fn=<DivBackward0>)
75 time traveller with a slight accession of cheerfulness really thi
76 Epoch -> 925 Loss -> tensor(0.0370, device='cuda:0', grad_fn=<DivBackward0>)
77 time traveller for so it will be convenient to speak of him was e
78 Epoch -> 950 Loss -> tensor(0.0311, device='cuda:0', grad_fn=<DivBackward0>)
79 time traveller for so it will be convenient to speak of him was e
80 Epoch -> 975 Loss -> tensor(0.0231, device='cuda:0', grad_fn=<DivBackward0>)
81 time traveller for so it will be convenient to speak of him was e
82 Epoch -> 1000 Loss -> tensor(0.0241, device='cuda:0', grad_fn=<DivBackward0>)
83 time traveller you can show black is white by argument said filby
84 Epoch -> 1025 Loss -> tensor(0.0303, device='cuda:0', grad_fn=<DivBackward0>)
85 time traveller for so it will be convenient to speak of him was e
86 Epoch -> 1050 Loss -> tensor(0.0304, device='cuda:0', grad_fn=<DivBackward0>)
87 time traveller for so it will be convenient to speak of him was e
88 Epoch -> 1075 Loss -> tensor(0.0278, device='cuda:0', grad_fn=<DivBackward0>)
89 time traveller you can show black is white by argument said filby
90 Epoch -> 1100 Loss -> tensor(0.0228, device='cuda:0', grad_fn=<DivBackward0>)
91 time traveller with a slight accession of cheerfulness really thi
92 Epoch -> 1125 Loss -> tensor(0.0267, device='cuda:0', grad_fn=<DivBackward0>)
93 time traveller you can show black is white by argument said filby
94 Epoch -> 1150 Loss -> tensor(0.0268, device='cuda:0', grad_fn=<DivBackward0>)
95 time traveller you can show black is white by argument said filby
96 Epoch -> 1175 Loss -> tensor(0.0173, device='cuda:0', grad_fn=<DivBackward0>)
97 time traveller with a slight accession of cheerfulness really thi
98 Epoch -> 1200 Loss -> tensor(0.0271, device='cuda:0', grad_fn=<DivBackward0>)
99 time traveller you can show black is white by argument said filby
100 Epoch -> 1225 Loss -> tensor(0.0307, device='cuda:0', grad_fn=<DivBackward0>)
101 time traveller for so it will be convenient to speak of him was e
102 Epoch -> 1250 Loss -> tensor(0.0257, device='cuda:0', grad_fn=<DivBackward0>)
103 time traveller fithed in familh imuthicesorlshith tanepsocr f ac
104 Epoch -> 1275 Loss -> tensor(0.0283, device='cuda:0', grad_fn=<DivBackward0>)
105 time traveller with a slight accession of cheerfulness really thi
106 Epoch -> 1300 Loss -> tensor(0.0268, device='cuda:0', grad_fn=<DivBackward0>)
```

```

10 time traveller you can show black is white by argument said filby
10 Epoch -> 1325 Loss -> tensor(0.0259, device='cuda:0', grad_fn=<DivBackward0>)
10 time traveller with a slight accession of cheerfulness really thi
19 Epoch -> 1350 Loss -> tensor(0.0239, device='cuda:0', grad_fn=<DivBackward0>)
10 time traveller with a slight accession of cheerfulness really thi
11 Epoch -> 1375 Loss -> tensor(0.0245, device='cuda:0', grad_fn=<DivBackward0>)
12 time traveller with a slight accession of cheerfulness really thi
13 Epoch -> 1400 Loss -> tensor(0.0299, device='cuda:0', grad_fn=<DivBackward0>)
14 time traveller with a slight accession of cheerfulness really thi
15 Epoch -> 1425 Loss -> tensor(0.0240, device='cuda:0', grad_fn=<DivBackward0>)
16 time traveller with a slight accession of cheerfulness really thi
17 Epoch -> 1450 Loss -> tensor(0.0226, device='cuda:0', grad_fn=<DivBackward0>)
18 time traveller for so it will be convenient to speak of him was e
19 Epoch -> 1475 Loss -> tensor(0.0278, device='cuda:0', grad_fn=<DivBackward0>)
10 time traveller for so it will be convenient to speak of him was e
12 Epoch -> 1500 Loss -> tensor(0.0240, device='cuda:0', grad_fn=<DivBackward0>)
12 time traveller for so it will be convenient to speak of him was e
13 Epoch -> 1525 Loss -> tensor(0.0227, device='cuda:0', grad_fn=<DivBackward0>)
12 time traveller with a slight accession of cheerfulness really thi
19 Epoch -> 1550 Loss -> tensor(0.0164, device='cuda:0', grad_fn=<DivBackward0>)
10 time traveller for so it will be convenient to speak of him was e
12 Epoch -> 1575 Loss -> tensor(0.0199, device='cuda:0', grad_fn=<DivBackward0>)
18 time traveller for so it will be convenient to speak of him was e
19 Epoch -> 1600 Loss -> tensor(0.0245, device='cuda:0', grad_fn=<DivBackward0>)
10 time traveller for so it will be convenient to speak of him was e
13 Epoch -> 1625 Loss -> tensor(0.0233, device='cuda:0', grad_fn=<DivBackward0>)
13 time traveller with a slight accession of cheerfulness really thi
13 Epoch -> 1650 Loss -> tensor(0.0213, device='cuda:0', grad_fn=<DivBackward0>)
13 time traveller with a slight accession of cheerfulness really thi
13 Epoch -> 1675 Loss -> tensor(0.0202, device='cuda:0', grad_fn=<DivBackward0>)
10 time traveller with a slight accession of cheerfulness really thi
13 Epoch -> 1700 Loss -> tensor(0.0203, device='cuda:0', grad_fn=<DivBackward0>)
18 time traveller you can show black is white by argument said filby
19 Epoch -> 1725 Loss -> tensor(0.0202, device='cuda:0', grad_fn=<DivBackward0>)
10 time traveller you can show black is white by argument said filby
14 Epoch -> 1750 Loss -> tensor(0.0240, device='cuda:0', grad_fn=<DivBackward0>)
12 time traveller with a slight accession of cheerfulness really thi
13 Epoch -> 1775 Loss -> tensor(0.0230, device='cuda:0', grad_fn=<DivBackward0>)
14 time traveller for so it will be convenient to speak of him was e
15 Epoch -> 1800 Loss -> tensor(0.0178, device='cuda:0', grad_fn=<DivBackward0>)
16 time traveller you can show black is white by argument said filby
17 Epoch -> 1825 Loss -> tensor(0.0249, device='cuda:0', grad_fn=<DivBackward0>)
18 time traveller you can show black is white by argument said filby
19 Epoch -> 1850 Loss -> tensor(0.0219, device='cuda:0', grad_fn=<DivBackward0>)
10 time traveller with a slight accession of cheerfulness really thi
15 Epoch -> 1875 Loss -> tensor(0.0207, device='cuda:0', grad_fn=<DivBackward0>)
13 time traveller you can show black is white by argument said filby
13 Epoch -> 1900 Loss -> tensor(0.0169, device='cuda:0', grad_fn=<DivBackward0>)
15 time traveller fir so it will be one four of emit do cane t whe
15 Epoch -> 1925 Loss -> tensor(0.0191, device='cuda:0', grad_fn=<DivBackward0>)
16 time traveller for so it will be convenient to speak of him was e
13 Epoch -> 1950 Loss -> tensor(0.0260, device='cuda:0', grad_fn=<DivBackward0>)
18 time traveller with a slight accession of cheerfulness really thi
10 Epoch -> 1975 Loss -> tensor(0.0215, device='cuda:0', grad_fn=<DivBackward0>)
10 time traveller with a slight accession of cheerfulness really thi
16 Epoch -> 2000 Loss -> tensor(0.0210, device='cuda:0', grad_fn=<DivBackward0>)
18 perplexity -> 1.0212559638748901 -*- 499128.1968430898 tokens/sec on cuda
18 time traveller with a slight accession of cheerfulness really thi
16 traveller with a slight accession of cheerfulness really this is

```

