## 10.6.  Recurrent neural networks

Suppose we're given a dataset of strings, and we want to be able to generate new strings of the same general type.

| | | |
|---|---|---|
| aaron | abraham | admatha |
| abdiel | achaia | adria |
| abel | adam | ahab |
| abigail | adlai | . . . |

Let's write a string $x$ as a sequence of characters, $x = x_1 x_2 \cdots x_n$. We want to be able to generate random sequences $X = X_1 X_2 \cdots X_N$ where both the characters and the length are random. We can take exactly the same approach we've taken to every single probability modelling problem throughout this course: first invent a parameterized distribution $\Pr_X(x\,;\,\theta)$, then fit the model, then use the fitted model to generate new random sequences. The quality of the results will of course depend on the model we choose. No model is true—but some models fit the data better than others, and we can evaluate the fit in the usual way, by the log likelihood of holdout data.

### OTHER PROBABILITY MODELS FOR SEQUENCES

The recurrent neural network can be thought of as a probability model for generating random sequences, designed to address shortcomings of other sequence models. It's useful to first review those other models. The simplest sequence model is a Markov chain.

Let's augment the alphabet with $\varnothing$ to denote 'start of string' and $\square$ to denote 'end of string', and generate random strings by starting at $\varnothing$ and jumping from character to character until we hit $\square$. The causal diagram is

$$\varnothing \to X_1 \to X_2 \to \cdots \to X_N \to \square$$

and the likelihood function is

$$\Pr_X(x_1 x_2 \cdots x_n\,;\,\theta) = p_\theta(x_1 \mid \varnothing)\, p_\theta(x_2 \mid x_1) \times \cdots \times p_\theta(\square \mid x_n)$$

where $\theta$ governs the character-to-character transition probabilities.

We'll get better results if we include some history, as per Markov's trigram model.

The causal diagram for Markov's trigram model, example 10.2.2, is
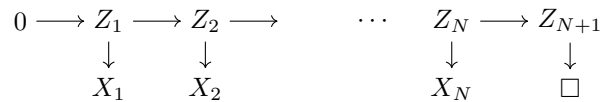
$$
\begin{array}{ccccccc}
\varnothing\varnothing \dashrightarrow & \varnothing X_1 \dashrightarrow & & \cdots & X_{N-2}X_{N-1} \dashrightarrow & X_{N-1}X_N & \\
\downarrow & \downarrow & & & \downarrow & \downarrow & \\
X_1 & X_2 & & & X_N & \square &
\end{array}
$$

and the likelihood function is

$$\Pr_X(x_1 x_2 \cdots x_n\,;\,\theta) = p_\theta(x_1 \mid \varnothing\varnothing)p_\theta(x_2 \mid \varnothing x_1) \times \cdots \times p_\theta(\square \mid x_{n-1}x_n).$$

In this diagram the solid arrows denote random generation of new values, and the dotted arrows denote inputs to the deterministic function that updates the 'last two characters' state variable.

How much history should we keep? The more we keep, the more faithful our random sequences are to the training dataset—but if we're too faithful then we'll just regurgitate the training data rather than generate novel sequences, so we'll assign a low likelihood to the novel datapoints in the holdout set, so we'll score badly on evaluation. Rather than hard-coding a fixed history window, a flexible alternative is to use a hidden Markov model, where the hidden state is some abstract summary of everything that's worth remembering about the past.

$$0 \longrightarrow Z_1 \longrightarrow Z_2 \longrightarrow \qquad \cdots \qquad Z_N \longrightarrow Z_{N+1}$$
$$\downarrow \qquad \downarrow \qquad\qquad\qquad \downarrow \qquad \downarrow$$
$$X_1 \qquad X_2 \qquad\qquad\qquad X_N \qquad \square$$

In this particular version we've included an initial hidden state $0$, to allow the first productive hidden state $Z_1$ to be random. The likelihood function is
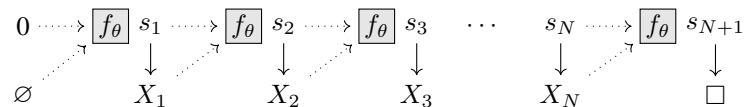
$$\mathrm{Pr}_X(x_1 x_2 \cdots x_n \; ; \; \theta) = \sum_{z_1,\ldots,z_{n+1}} \left\{ \begin{array}{l} p_\theta(z_1|0)\, q_\theta(x_1|z_1) \times p_\theta(z_2|z_1)\, q_\theta(x_2|z_2) \times \cdots \\ \cdots \times p_\theta(z_{n+1}|z_n)\, q_\theta(\square|z_{n+1}). \end{array} \right\}$$

where $p$ is the transition probabilities for the hidden state, and $q$ is the emission probabilities for reading off characters, and $\theta$ represents all the parameters that need to be estimated for both these functions.

On one hand the hidden Markov is wonderfully flexible because we can choose the hidden state space however we like. On the other hand this is a useless idea because it doesn't solve anything, it just passes the buck on to whoever has to design the hidden state transitions; and because we've ended up with a horrible intractable sum for the likelihood function. This probability model will not be easy to work with when it comes to fitting the distribution.

## RECURRENT NEURAL NETWORKS TO DIGEST HISTORY

It's useful to be able to keep track of the process's history, but it's hard to know exactly how history should be represented. What if we could train a neural network to learn this? That's exactly the idea behind the recurrent neural network probability model. It has the same causal diagram structure as the trigram model, but instead of a hard-coded history window it allows a flexible representation of history via a neural network and a 'history digest' $s_i$ that can have as many dimensions as we like.

$$0 \dashrightarrow \boxed{f_\theta} \; s_1 \dashrightarrow \boxed{f_\theta} \; s_2 \dashrightarrow \boxed{f_\theta} \; s_3 \quad \cdots \quad s_N \dashrightarrow \boxed{f_\theta} \; s_{N+1}$$
$$\downarrow \qquad\quad \downarrow \qquad\quad \downarrow \qquad\qquad \downarrow \qquad\qquad \downarrow$$
$$\varnothing \qquad\quad X_1 \qquad\quad X_2 \qquad\quad X_3 \qquad\quad X_N \qquad\quad \square$$

Here, $f_\theta$ is some neural network that takes in the latest event $X_{i-1}$ and the previous history digest $s_{i-1}$, and computes a new history digest $s_i$. The next event $X_i$ is then generated randomly, using $s_i$ as its parameters. To save the bother of having to specify how $X_i$ is to be generated, we might as well split $s_i$ into two parts, $s_i = (p_i, v_i)$ where $p_i$ specifies the likelihood of $X_i$ and $v_i$ is used for computing $s_{i+1}$. This leads to the basic recurrence equations

$$(p_i, v_i) = f_\theta(x_{i-1}, v_{i-1}), \qquad X_i \sim \mathrm{Cat}(p_i).$$

Generating random sequences.   Once we've trained the neural network $f_\theta$ it's easy to generate a random string.

```
def generate():
    x,v = "∅",0
    while x[−1] ≠ □:
        p,v = fθ(x[−1], v)
        newchar = np.random.choice(ALPHABET, p)
        x = x + newchar
    return x[1:−1]   # strip out ∅ and □
```

Training the network.   The neural network is easily trained because, like the trigram model, there is an explicit formula for the log likelihood of a datapoint. The key is that the states are visible, not hidden—they're patent, not latent—in the sense that given a datapoint $x = x_1 x_2 \cdots x_n$ we can compute $s_1, \ldots, s_n$. In pseudocode the log likelihood function looks like this:

```
def loglik(x):
    res = 0
    lastchar,v = ∅,0
    for char in x + "□":
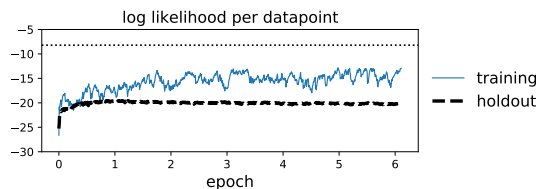```

```
        p,v = f_θ(lastchar, v)
        res = res + log(p[char])
        lastchar = char
    return res
```

Evaluating the model.    The perfectly-fitted probability model for the training dataset is its empirical distribution, which assigns probability $1/N$ to every one of the $N$ datapoints in the whole dataset (assuming no duplication). The best possible log likelihood we can possibly achieve during training is thus $N \log(1/N)$. This can give us a hint about whether our network is overfitting, or whether it's too simple to overfit:



- If the training loss gets close to this bound, it suggests the network is overfitted to the training data, and so the log likelihood of the holdout dataset will be low. This suggests we need to add a regularizing mechanism such as dropout.
- If the neural network can't get close to this even with the regularizing mechanism turned off, it suggests the network isn't complex enough to describe the dataset, so we need more nodes or layers.

$$* \; \maltese \; *$$

Here is PyTorch code, a bit more fussy than our pseudocode. The convention for PyTorch sequence operations is to act on tensors of shape $n \times b \times f$ where we're working on a batch of $b$ sequences at a time, $n$ is the maximum length of these sequences, and $f$ is the dimensionality of each item in a sequence. Also, PyTorch manages the iteration for us: we just call $f_\theta(x_0 x_1 \cdots x_{n-1}, v_0)$ and it doesn't just return $(p_1, v_1) = f_\theta(x_0, v_0)$ it returns $((p_1, \ldots, p_n), v_n)$.

```
1   url = "https://www.cl.cam.ac.uk/teaching/2021/DataSci/data/english_names.txt"
2   names = pandas.read_table(url, header=None, names=['X']).X.str.lower()
3   alphabet = list(set(''.join(names)))
4   alphabet_n = {c:i+1 for i,c in enumerate(alphabet)}
5   A = len(alphabet)+1
6
7   # Convert from strings to integer-coding (with 0 at beginning and end) then one-hot coding
8   names_n = [[0] + [alphabet_n[x] for x in name] + [0] for name in names]
9   names_n = [torch.tensor(n) for n in names_n]
10  names_oh = [nn.functional.one_hot(n, num_classes=A).float() for n in names_n]
11
12  class RSeq(nn.Module):
13      def __init__(self, H=50, L=2):
14          super().__init__()
15          # Use rnn for the iterative part, producing a vector of dim. H per timestep,
16          # then apply an additional map before turning it into a probability distribution of dim. A
17          self.rnn = nn.GRU(input_size=A, hidden_size=H, num_layers=L, dropout=.05)
18          self.map = nn.Linear(self.rnn.hidden_size, A)
19
20      def f(self, x, v=None):   # x.shape m × b × A, v.shape L × b × H
21          assert x.shape[1]==1, "This code only works with batch_size=1"
22          y,v = self.rnn(x, v) # y.shape m × b × H, v.shape L × b × H
23          z = self.map(y)        # z.shape m × b × A
24          logp = nn.functional.log_softmax(z, dim=2) # logp.shape m × b × A
25          return logp,v
26
27      def forward(self, x): # x.shape = (n+1)*b*A, x assumed to have ∅ and □
28          logp,_ = self.f(x[:-1])        # logp.shape n × b × A
29          return torch.sum(logp * x[1:], dim=(0,2)) # shape b
30
31      def generate(self):
32          is_training = self.training
```
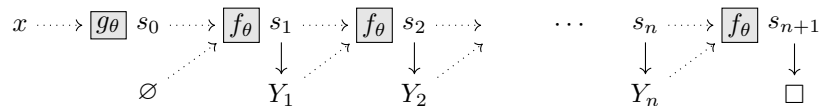
```
33              self.train(False)
34              with torch.no_grad():
35                  res = ''
36                  x_n,v = 0,None
37                  while True:
38                      x_n = torch.tensor(x_n).reshape(1,1)
39                      x_oh = nn.functional.one_hot(x_n, num_classes=A).float()
40                      logp,v = model.f(x_oh, v)
41                      p = torch.exp(logp[0,0])
42                      x_n = np.random.choice(A, p=p.detach().numpy())
43                      if x_n == 0: break
44                      res += alphabet[x_n−1]
45              self.train(is_training)
46              return res
47
48  model = RSeq()
49  model.train(True)
50  optimizer = optim.Adam(model.parameters())
51
52  with Interruptable() as check_interrupted:
53      while True:
54          for j in np.random.permutation(len(names_oh)):
55              check_interrupted()
56              optimizer.zero_grad()
57              e =− model(names_oh[j][:,None,:])
58              e.backward()
59              optimizer.step()
```

$$* \divideontimes *$$

Recurrent neural networks easily be used for supervised learning tasks. Suppose for example we want to train a neural network to write captions for images. Let the the image be $x$, model the caption as a random sequence $Y = Y_1 Y_2 \cdots Y_n$, and fit a model with an extra neural network $g_\theta$ for processing the image and converting it to an initial state $s_0$.



The Transformer[63] architecture, a neural network that generates text much more impressively than does a recurrent neural network, also uses sequential generation with non-hidden state, but with an even simpler design. It generates each item in the sequence one by one, using the model $X_n \sim \text{Cat}(f_\theta(x_1 x_2 \cdots x_{n-1}))$. The cleverness consists in designing a neural network that can accept an arbitrary length sequence as its input. But from the point of view of random variables, there's barely any difference between a recurrent neural network and a transformer.

---

[63]Ashish Vaswani et al. "Attention Is All You Need". In: *NIPS*. 2017. URL: http://arxiv.org/abs/1706.03762.