

RNN LANGUAGE MODELS

- ① We have a lot of predictive power, as humans, when it comes to language. Suppose you're asked to guess the next letter in the sequence:

W A _

Likely you don't have a uniform probability over letters. Letters like S or N are quite likely (in English), but Q or E are pretty unlikely (this, despite the fact that E is generally a common letter).

- ② More context can sharpen the probability:

G E O R G E _ W A _

B E E S W A _

W H A T _ D O _ Y O U _ W A _

- ③ A language model over an alphabet A is a probability distribution over letter $a_{n+1} \in A$ given the previous letters in a sequence a_1, a_2, \dots, a_n .

$$\Pr(a_{n+1} | a_1, a_2, \dots, a_n)$$

RNN LANGUAGE MODELS

- ④ Language models are frequently used in natural language processing to determine how "fluent" a sequence sounds (relative to a particular language). For instance, a good language model (LM) for English should have

$$\Pr("N"|"W", "A") > \Pr("E"|"W", "A")$$

- ⑤ A sequence's probability can be scored as:

$$\Pr(a_1, \dots, a_k) = \prod_{k=1}^k \Pr(a_k | a_1, \dots, a_{k-1})$$

A good LM for English should have

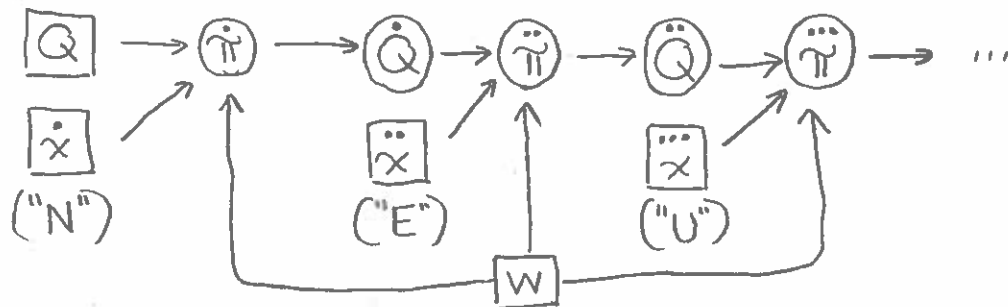
$$\Pr("THERE IS A") \gg \Pr("IHRA S EET")$$

- ⑥ RNNs provide an effective way to create LMs.

The main difference between RNN LMs and the RNNs we've seen so far is that RNN LMs don't just take an input character at each time step, they also emit an output character.

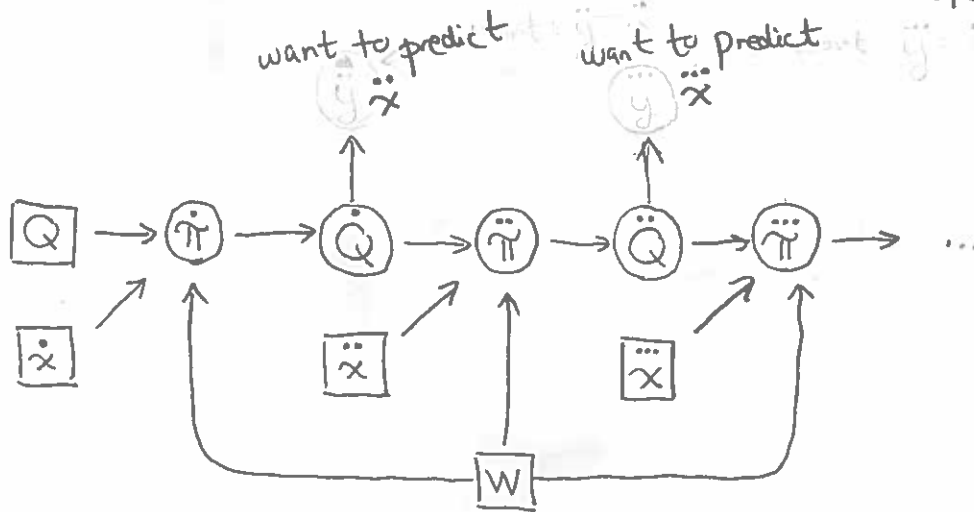
RNN LANGUAGE MODELS

⑦ Consider our standard RNN:



As it reads in letters ("NEU..."), it updates the state vector Q .

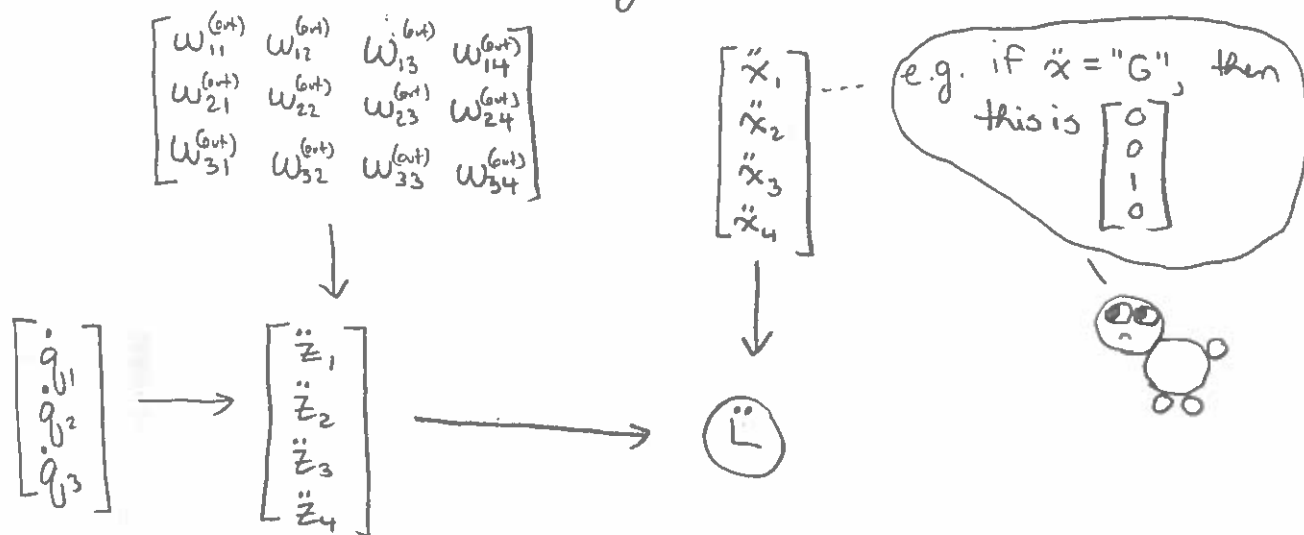
⑧ To make this into a language model, we will try to predict the next letter from each state vector:



RNN LANGUAGE MODELS

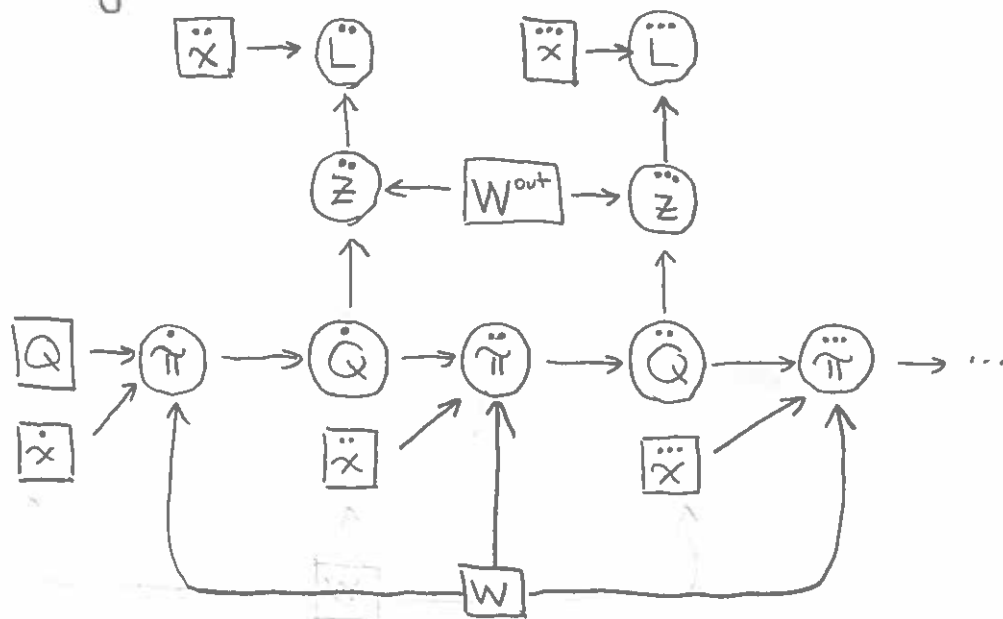
⑨ To show some simple concrete examples, let's assume we're learning a language model over DNA sequences (so our alphabet is $\{A, C, G, T\}$).

To predict the next letter from the current state vector \vec{q} , we can use a simple dense layer:



where: $L(\vec{z}, \vec{x}) = -\log(\vec{x}^T \cdot \text{softmax}(\vec{z}))$
 $\vec{z} = W^T \vec{q}$

⑩ Putting this into our RNN:



RNN LANGUAGE MODELS

- ① It's mildly awkward that our loss is distributed across each layer, but it's easy to reconcile using probability and logs. Essentially, we want to maximize the predicted probability of the actual sequence:

$$\prod_k \Pr(a_k)$$
$$= \prod_k \mathbf{x}^{(k)T} \cdot \text{softmax}\left(\frac{\mathbf{b}}{Z}\right)$$

This is equivalent to maximizing the log probability

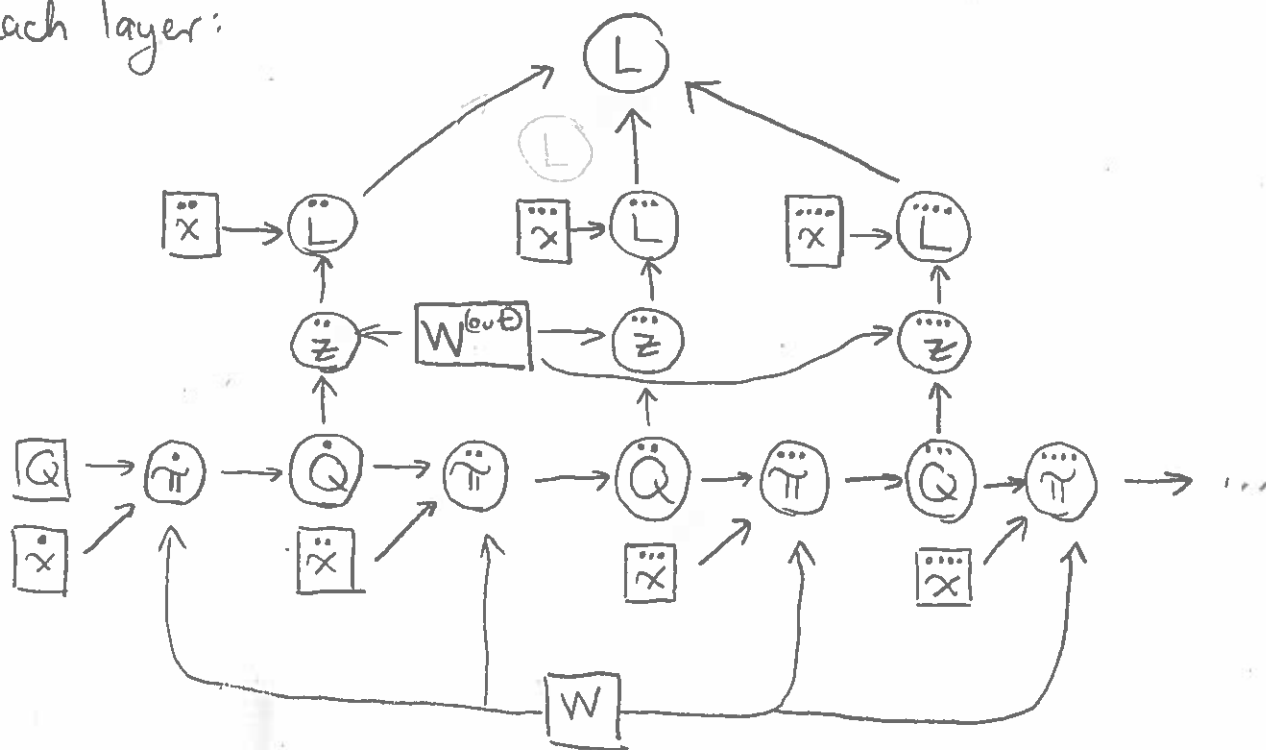
$$\log \prod_k \mathbf{x}^{(k)T} \cdot \text{softmax}\left(\frac{\mathbf{b}}{Z}\right)$$
$$= \sum_k \log \mathbf{x}^{(k)T} \cdot \text{softmax}\left(\frac{\mathbf{b}}{Z}\right)$$

Or alternatively, minimizing the negative log probability:

$$-\log \prod_k \mathbf{x}^{(k)T} \cdot \text{softmax}\left(\frac{\mathbf{b}}{Z}\right)$$
$$= -\sum_k \log \mathbf{x}^{(k)T} \cdot \text{softmax}\left(\frac{\mathbf{b}}{Z}\right)$$
$$= \sum_k -\log \mathbf{x}^{(k)T} \cdot \text{softmax}\left(\frac{\mathbf{b}}{Z}\right)$$
$$= \sum_k L\left(\frac{\mathbf{b}}{Z}, \mathbf{x}\right)$$

RNN LANGUAGE MODELS

12) So the overall loss is just the sum of the losses at each layer:



13) From this depiction, it is straightforward (though arduous by hand) to compute $\frac{\partial L}{\partial W}$ and $\frac{\partial L}{\partial W^{(out)}}$ for each sequence $\dot{x}, \ddot{x}, \ddot{\ddot{x}}, \dots$ in our training data.

$$\text{e.g. } \frac{\partial L}{\partial W} = \sum_{k=2}^K \frac{\partial L}{\partial L^{(k)}} \cdot \frac{\partial L^{(k)}}{\partial W}$$

= ... etc.