

Objectives:

- Formulate, estimate and sample sequence from Markov models.
- Understand the relation between RNNs and Markov model for generating sequences
- Understand the process of decoding of RNN in generating sequences.

Markov Models:

The k th order Markov model is a simple probability model of giving a distribution of what comes next, based on looking at k steps backwards.

Let $w \in V$ denote the set of possible words/symbols that includes:

- on **UNK** symbol for any unknown word (out of vocabulary)
- **<beg>** symbol for specifying the start of a sentence
- **<end>** symbol for specifying the end of the sentence

<beg> The lecture leaves me **UNK** **<end>**
 $w_0 \quad w_1 \quad w_2 \quad w_3 \quad w_4 \quad w_5 \quad w_6$

In a first order Markov model (bigram model) the next symbol only depends on the previous one.

$$P(w_1, \dots, w_6) = P(\underbrace{w_1}_{<beg>} | \underbrace{w_0}_{<beg>}) \dots P(\underbrace{w_6}_{<end>} | \underbrace{w_5}_{<end>})$$

	ML	course	is	UNK	<end>	
<beg>	0.7	0.1	0.1	0.1	0.0	$\rightarrow \sum = 1$
ML	0.1	0.5	0.2	0.1	0.1	
course	0.0	0.0	0.7	0.1	0.2	
is	0.1	0.3	0.0	0.6	0.0	
UNK	0.1	0.2	0.2	0.3	0.2	

<beg> course is ^{UNK} ~~great~~ <end>

$$P(\text{course} | <beg>) \quad P(\text{is} | \text{course}) \quad P(\text{UNK} | \text{is}) \quad P(<end> | \text{UNK})$$

0.1 0.7 0.6 0.2

<beg> course is UNK <end>

we must look at the row here, probability of distribution over the first word which comes after the beginning symbol.

To sample from it, we throw a weighted 5 way die, weighted heavily towards ML, a little bit less for course and is, and the end symbol has probability 0. Most likely we would get the word ML, but we might also get the word course.

- The goal is to maximize the probability that the model can generate all the observed sentences (corpus S)

$$s \in S, s = \{w_1^s, w_2^s, \dots, w_{|s|}^s\}$$

$$\prod_{s \in S} \left[\prod_{i=1}^{|s|} p(w_i^s | w_{i-1}^s) \right]$$

for all the sentences 1 sentence.

Maximum likelihood estimation.

- The ML estimate is obtained as normalized counts of successive word occurrences (matching statistics)

$$\text{count}(w, w')$$

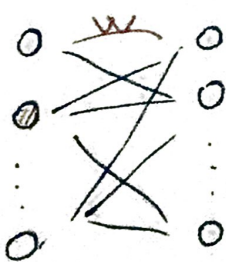
$$\hat{p}(w' | w) = \frac{\text{count}(w, w')}{\sum_{\tilde{w}} \text{count}(w, \tilde{w})}$$

Markov Models to Feed forward Neural Nets:

we can also represent the Markov model as a feed forward neural network (very extensible)

one hot feature vector $\vec{x}(w_{i-1})$

one hot vector
All 0 and 1 non-zero



$$p_k = P(w_i = k | w_{i-1}) \quad p_k \geq 0 \text{ and } \sum_k p_k = 1$$

↑
prob of next word

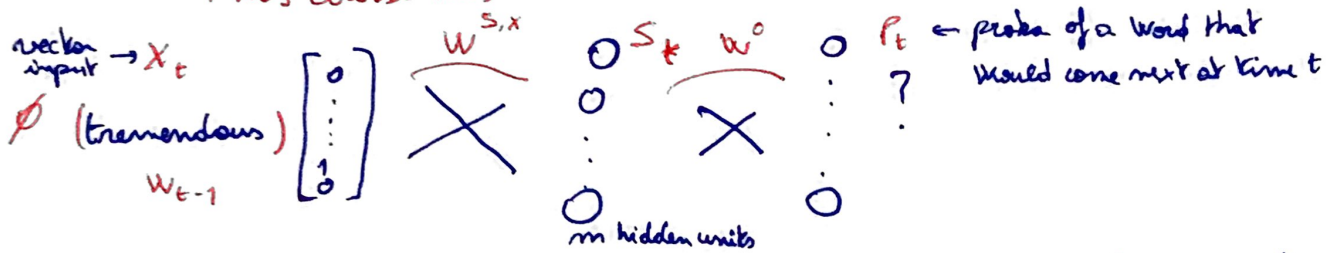
$$z_k = \sum_j x_j w_{jk} + w_{0k} \in \mathbb{R}$$

Softmax output layer: $P_k = \frac{e^{z_k}}{\sum_j e^{z_j}} \Rightarrow P_k \geq 0, \sum_k P_k = 1$

RNN Deeper Dive:

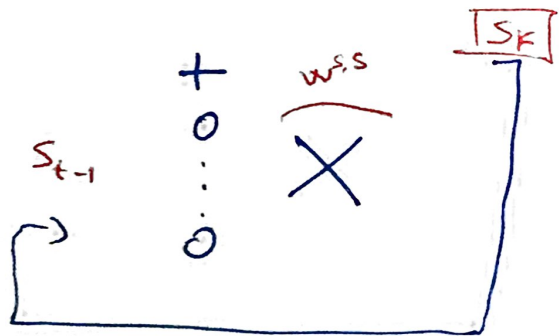
language modeling: What comes next?

This course has been a tremendous | ...



We feed into the model a one hot vector of the previous word, where t now refers to the time or position in the sentence when we are making a prediction.

We make this into a RNN model by feeding in the previous hidden state activity



$$s_t = \tanh(W^{s,s}s_{t-1} + W^{s,x}x_t) \text{ state}$$

$$P_t = \text{softmax}(W^o s_t) \text{ output distribution}$$

Decoding:

