

&lt; Notes



## Neural Network based LMs

Char-RNN

Announcement: Parts 4+5 of HW5 will now be part of HW6.  
 Parts 1-3 of HW5 are due on Wednesday.

Goal of language modeling:

- Learn a function that returns the joint probability of a sequence of words

Primary difficulties:

- ① There are too many parameters to accurately estimate.  
 This is sometimes called the "curse of dimensionality"
- ② In n-gram based models, we fail to generalize to related words / word sequences that we have observed.

### ① Curse of dimensionality / sparse statistics

Suppose we want a joint distribution over 10 words  
 " have a vocab of size 100,000

$$100,000^{10} = 10^{50} \text{ parameters}$$

too high to estimate from data

In LMs we use chain rule to get the conditional probability of next word in sequence given all of the previous words

$$p(w_1, w_2, w_3, \dots, w_T) = \prod_{t=1}^T p(w_t | w_1, \dots, w_{t-1})$$

What assumptions do we make in unigram LMs to simplify this? Probability of the next word only depends on previous N-1 words ← get an estimate



We construct tables to look up the



Notes

to simplify this? Probability of the next word only depends on previous  $N-1$  words get an estimate for small  $N$  from data

We construct tables to look up the probability of seeing a word given a history

history ( $N-1$ )	( curse of dimensionality )	$P(w_T   w_{T-n} \dots w_{T-1})$
$w_T$	azure	—
	oak	—
	knowledge	—

"STUPID" multinomial distributions == Table lookup.

What happens when we have a new (unseen) combination of  $N$  words, the table only stores observed sequences

- ① Back-off
- ② Smoothing / interpolation

We are basically just stitching together short sequences of observed words.

### Alternate idea

Let's try generalizing.

Intuition: take a sentence like

The cat is walking in the bedroom  
and use it when we assign probabilities to a similar sentence like

$\text{sim}(\text{cat}, \text{dog})$  The dog is running around a room.



&lt; Notes



Bengio et al 2003 NIPS paper "A Neural Probabilistic Language model".

$$M = 30, 60, 100$$

- prob. |
- LM
- ① use a vector space model where words are vectors w/ real values  $\mathbb{R}^m$ . Gives them a way to compute word similarity
  - ② Define a fn that returns a joint probability of words in a sequence based on a sequence of these vectors
  - ③ Simultaneously learn the word representations and the probability fn from data.

Seeing one of the cat/dog sentences allows them to increase the probability for that sentence and its combinatorial # of "neighbor" sentences in vector space.

A neural network architecture for prob. LMS.

Given: a training set  $w_1 \dots w_T$  where  
 $w_i \in V$

Learn:  $f(w_1 \dots w_T) = \hat{P}(w_T | w_1 \dots w_{T-1})$

s.t. giving a high probability to an unseen test/dev set (e.g. minimize the perplexity)

constraint: create a proper probability distribution (e.g. sums to 1) so that we can take the product of conditional probs to get the joint prob of a sentence.

Model

- ① Create a mapping function  $C$  from any word in  $V$  onto  $\mathbb{R}^m$ .  $m = 30, 60, 100$



→ Store this in a  $N$ -by- $M$  matrix



&lt; Notes



## Model

① Create a mapping function  $C$  from any word in  $V$  onto  $\mathbb{R}^M$ .  $M = 30, 60, 100$

Store this in a  $V$ -by- $M$  Matrix.

Initialize it w/SVD

② The neural architecture:

a fn  $g$  maps seq of word vectors onto a prob. distribution over words vocab  $V$ .

$$g(C(w_{+N}) \dots C(w_{+1})) = \hat{P}(w_t | w_{+N}, \dots, w_{+1})$$

$g$  is implemented w/ a feed forward or recurrent neural network.

trained w/stochastic gradient descent

they use "softmax" in output layer

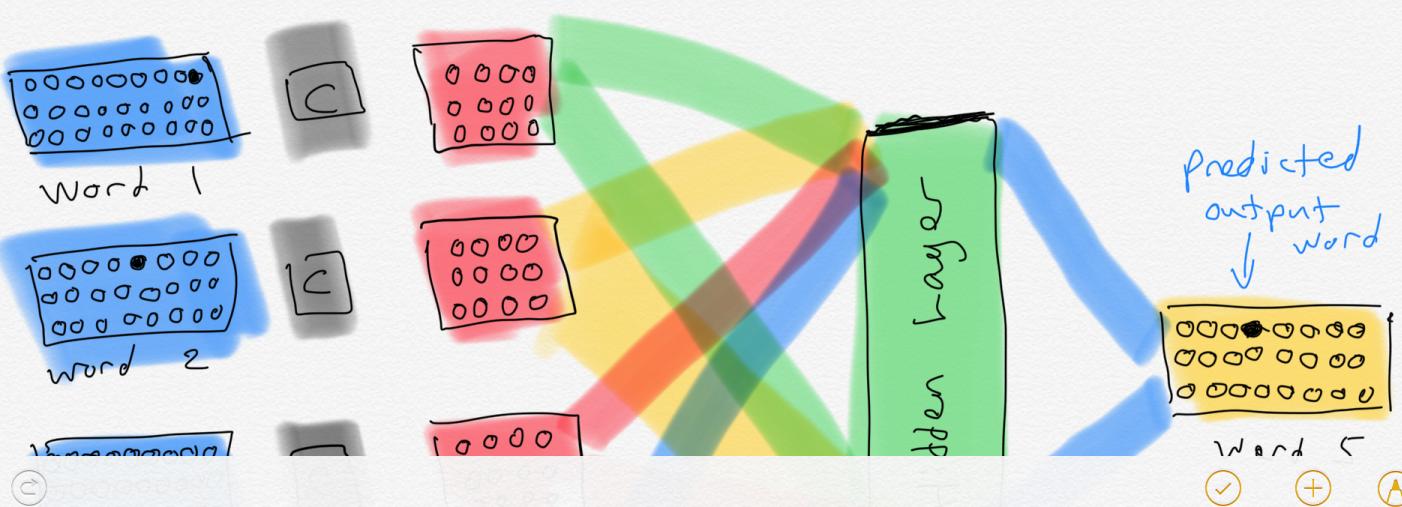
Claim: NN LM is better than the n-gram LM

Verify. = Evaluate

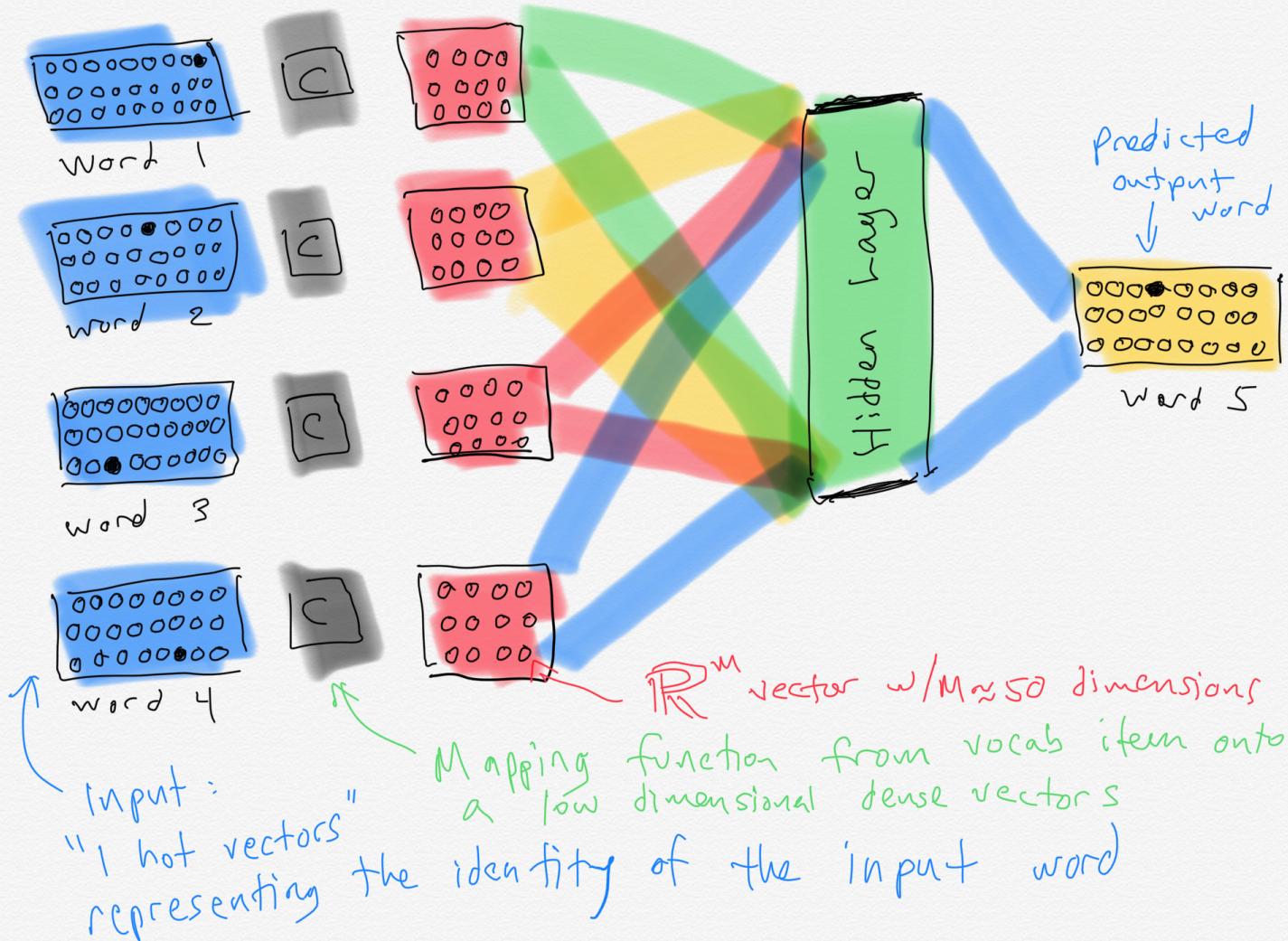
- Train on same training data

- Intrinsic evaluation or extrinsic eval

Perplexity on previously unseen data (test set)



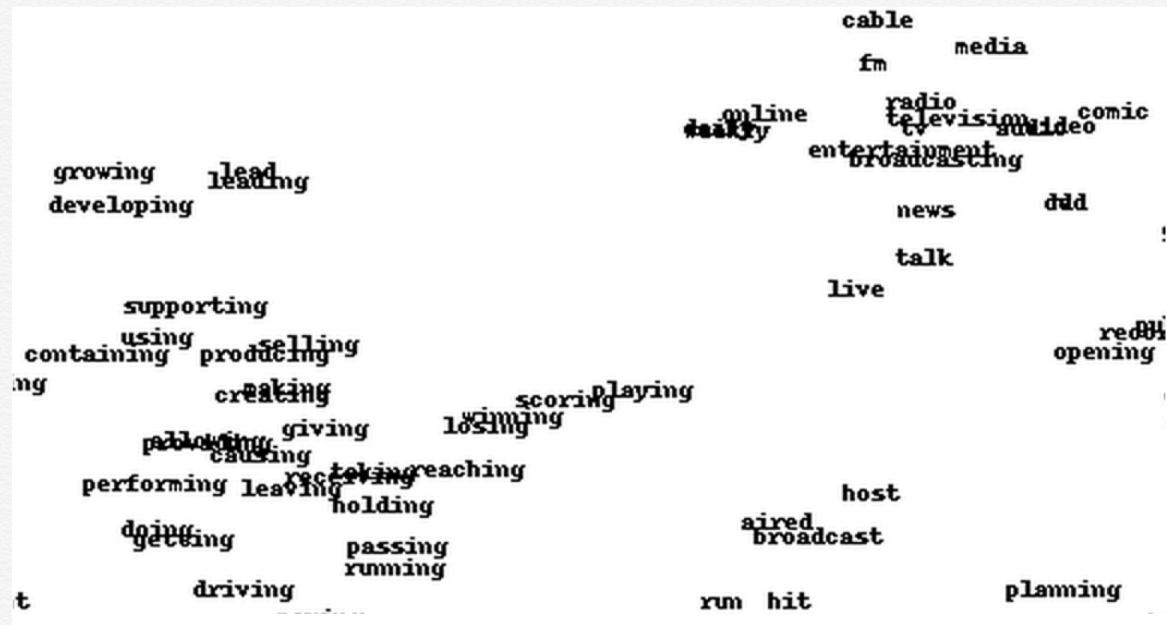
&lt; Notes



A by-product of training a neural network based language models are dense word embeddings. This is the way that word embeddings that you download for the advanced vector space models like were created. E.g. Word2vec, GLOVE.



&lt; Notes



When the ~50 dimensional vectors that result from training a neural LM are projected down to 2-dimensions, we see a lot of words that are intuitively similar to each other are close together.

For more information about Neural Networks, check out chapter 8 of the textbook.

It covers

- Neural network training
- activation functions
- feed forward NNs
- Training w/ loss funcs and stochastic gradient descent.