

# Natural Language Processing: High-Dimensional Outputs

Sargur N. Srihari  
srihari@cedar.buffalo.edu

This is part of lecture slides on [Deep Learning](http://www.cedar.buffalo.edu/~srihari/CSE676):  
<http://www.cedar.buffalo.edu/~srihari/CSE676>

# Topics in NLP

1. N-gram Models
2. Neural Language Models
3. High-Dimensional Outputs
4. Combining Neural Language Models with n-grams
5. Neural Machine Translation
6. Historical Perspective

# Topics in High-Dimensional Outputs

- Overview
  1. Use of a Short List
  2. Hierarchical Softmax
  3. Importance Sampling
  4. Noise-Contrastive Estimation and Ranking Loss

# Word vocabularies can be large

- In many NLP applications we want our models to produce words (rather than characters) as the fundamental output
- For large vocabularies it is computationally expensive to represent output distribution over the choice of a word because the vocabulary size is large
  - Ex: In many applications  $V$  contains 100K words

# Naiive Approach

- Naïve approach
  1. Apply an affine transformation from hidden representation to output space
  2. Then apply a softmax function from hidden representation to output space
- The Weight matrix representing this affine transformation is very large because output dimension is  $|V|$ 
  - High memory cost to represent it and high computational cost to multiply by it

# Cost at both training and testing

- Because softmax is normalized across all  $|V|$  outputs, it is necessary to perform full matrix multiplication at training time and testing time
- We cannot calculate only the dot product with the weight vector for the correct output
- Thus the high computational cost of the output layer arises at both
  - Training time (to compute likelihood and gradient)
  - Testing time (to compute probabilities for all selected words)

# Cost of Naïve Mapping to Words

- Suppose  $\mathbf{h}$  is the top hidden layer used to predict output probabilities  $\hat{y}_i$
- If we parameterize the transformation from  $\mathbf{h}$  to  $\hat{y}_i$  with learned weights  $W$  and learned biases  $\mathbf{b}$ , then the affine-softmax layer performs the following computations

$$\begin{aligned} a_i &= b_i + \sum_j W_{ij} h_j \quad \forall i \in \{1, \dots, |V|\} \\ \hat{y}_i &= \frac{e^{a_i}}{\sum_{i'=1}^{|V|} e^{a_{i'}}} \end{aligned}$$

- If  $\mathbf{h}$  contains  $n_h$  elements then the above operation is  $O(|V|n_h)$ 
  - $n_h$  is in the thousands and  $|V|$  is in hundreds of thousands

# Use of a Short List

- First neural language models to deal with high cost of using softmax over large  $V$ :
  - Split  $V$  into shortlist  $L$  of frequent words (say 10,000 words handled by a neural net) and tail  $T=V/L$  of rare words (handled by an n-gram model)
  - To combine two predictions, NN has to predict probability that word appearing after context  $C$  belongs to tail list
    - By extra sigmoid output unit to provide an estimate
    - The extra output can then be used to achieve an  $\text{est}P(i \in T | C)$ . probability estimate over all words in  $V$  as follows:
 

$$P(y = i | C) = 1_{i \in L} P(y = i | C, i \in L) (1 - P(i \in T | C)) + 1_{i \in T} P(y = i | C, i \in T) P(i \in T | C)$$
    - where  $P(y=i|C, i \in L)$  is provided by neural language model
    - and  $P(y=i|C, i \in T)$  is provided by the n-gram model



# Disadvantage of Short-list approach

- Potential generalization advantage of the neural language models is limited to the most frequent words
  - Where it is least useful
- This disadvantage has stimulated exploration of alternative methods to deal with high-dimensional outputs
  - Described next

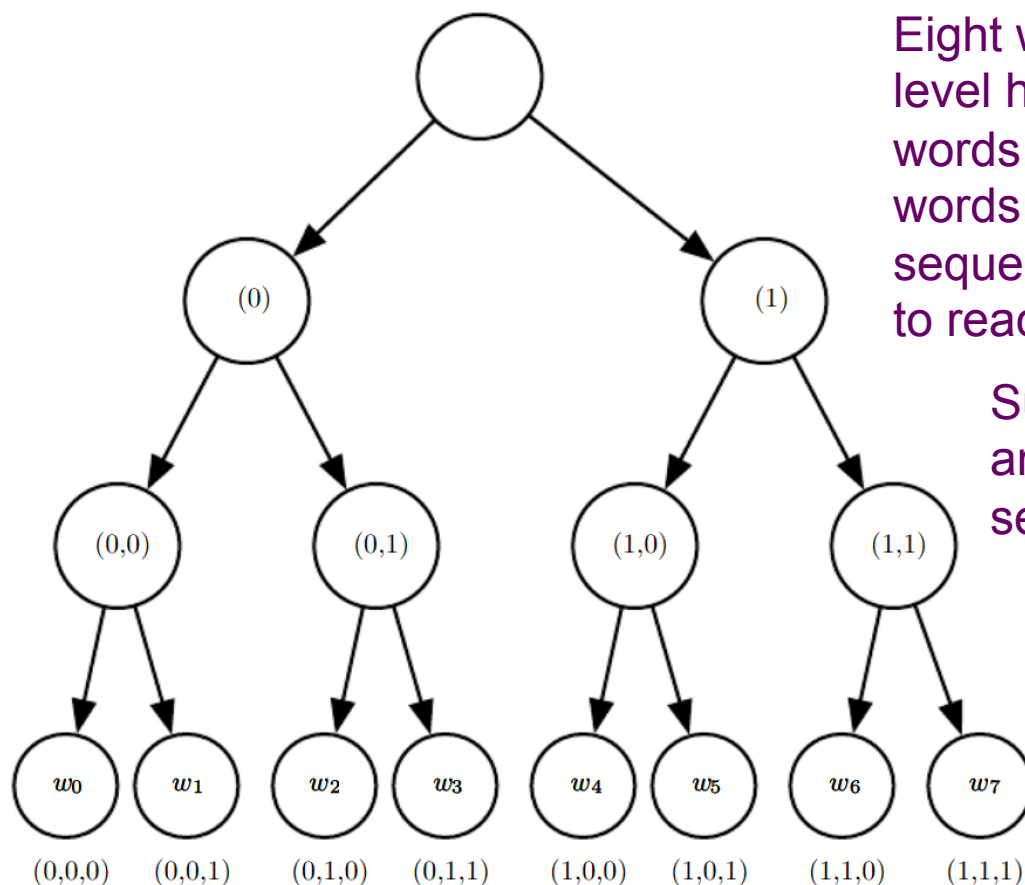
# Hierarchical Softmax

- Classical approach to reducing computational burden of high-dimensional outputs over large vocabulary sets  $V$  is to decompose probabilities hierarchically
- Instead of necessitating a no. of computations proportional to  $|V|$  (and also proportional to no. of computations proportional to no. of hidden units  $n_h$ ), the  $|V|$  factor can be reduced to as low as  $\log |V|$

# Hierarchy of Words

- Hierarchy builds categories of words
- Then categories of categories of words, etc
- Nested categories form a tree
  - With words at the leaves
- In a balanced tree, tree has depth  $O(\log|V|)$
- Probability of choosing a word is given by:
  - The product of the probabilities of choosing the branch leading to that word at every node on a path from the root of the tree to the leaf containing the word
  - A simple example is given next

# Simple Hierarchy of Word Categories



Eight words  $w_0, \dots, w_7$  organized into a three level hierarchy. Leaves represent specific words. Internal nodes represent groups of words. Any node can be indexed by the sequence of binary decisions (0=left, 1=right) to reach the node from the root

Superclass (0) contains the classes (0,0) and (0,1) which respectively contain the sets of words  $\{w_0, w_1\}$  and  $\{w_2, w_3\}$

Superclass (1) contains the classes (1,0) and (1,1) which respectively contain the sets of words  $\{w_4, w_5\}$  and  $\{w_6, w_7\}$

Node (1,0) corresponds to the prefix  $(b_0(w_4)=1, b_1(w_4)=0)$  and the probability of  $w_4$  can be decomposed as:

$$\begin{aligned}
 P(y = w_4) &= P(b_0 = 1, b_1 = 0, b_2 = 0) \\
 &= P(b_0 = 1)P(b_1 = 0 \mid b_0 = 1)P(b_2 = 0 \mid b_0 = 1, b_1 = 0)
 \end{aligned}$$

# Importance Sampling

- Training of neural language models can be speeded up by avoiding explicit computation of contribution to the gradient from all words that do not appear in the next position
- Every incorrect word should have low probability under the model
- Instead of enumerating all words, it is possible to sample only a subset of words

# Gradient using Sampling

Using notation:

$$a_i = b_i + \sum_j W_{ij} h_j \quad \forall i \in \{1, \dots, |\mathbb{V}|\},$$

$$\hat{y}_i = \frac{e^{a_i}}{\sum_{i'=1}^{|\mathbb{V}|} e^{a_{i'}}}.$$

where  $\mathbf{a}$  is the vector of pre-softmax activations (or scores) with one element per word.

the gradient written as follows:

$$\begin{aligned} \frac{\partial \log P(y | C)}{\partial \theta} &= \frac{\partial \log \text{softmax}_y(\mathbf{a})}{\partial \theta} \\ &= \frac{\partial}{\partial \theta} \log \frac{e^{a_y}}{\sum_i e^{a_i}} \\ &= \frac{\partial}{\partial \theta} (a_y - \log \sum_i e^{a_i}) \\ &= \frac{\partial a_y}{\partial \theta} - \sum_i P(y = i | C) \frac{\partial a_i}{\partial \theta} \end{aligned}$$

The first term is the **positive phase term**  
(pushing  $a_y$  up)

Second term is the **negative phase term**  
(pushing  $a_i$  down for all  $i$  with weight  $P(i | C)$ )

Since negative phase is expectation, can estimate with a Monte Carlo sample

# Sampling from another distribution

- Gradient method based on sampling would require sampling from the model itself
- Sampling from the model requires computing  $P(i|C)$  for all  $i$  in the vocabulary
  - Which is precisely what we are trying to avoid
- Instead of sampling from the model, one can sample from another distribution, called the proposal distribution (denoted  $q$ )
  - And use weights to correct for bias due to sampling from wrong distribution
  - This is an application of importance sampling

# Biased Importance Sampling

- Even exact importance sampling is inefficient
- It requires computing weights  $p_i/q_i$  where  $p_i = P(i|C)$ 
  - Which can only be computed if all the scores  $a_i$  are computed
- The solution adopted is called biased importance sampling
  - Where the importance weights are normalized to sum to 1