

Word Embedding using Word2Vec

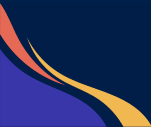
Ramaseshan Ramachandran

TABLE OF CONTENTS I

- ① Word2Vec
- ② One-Word Learning
 - CBOW Model for multiple words
 - Loss function
 - What does it learn?
 - Skip-Gram model
 - Sub-sampling

- Softmax
- Hierarchical Softmax
- Softmax
- Hierarchical Softmax - Architecture
- Negative Sampling
- Limitations of Word2Vec

- ③ References



GOAL

- ▶ Process each word in a Vocabulary of words to obtain a respective numeric representation of each word in the Vocabulary
- ▶ Reflect semantic similarities, Syntactic similarities, or both, between words they represent
- ▶ Map each of the plurality of words to a respective vector and output a single merged vector that is a combination of the respective vectors

CONTEXT WORDS AND CENTRAL WORD

Continuous Bag of Words (CBOW) Model – A central word is surrounded by context words. Given the context words identify the central word

- ▶ Wish you many

more	happy	returns	of	the
------	-------	---------	----	-----

 day

Skip Gram Model – Given the central word, identify the surrounding words

- ▶ Wish you many

more	happy	returns	of	the
------	-------	---------	----	-----

 day

CONTINUOUS BAG OF WORDS (CBOW)

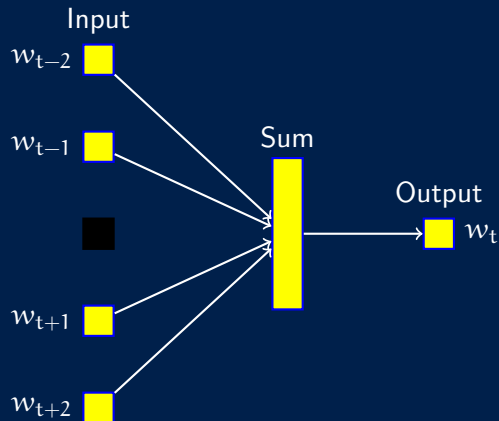
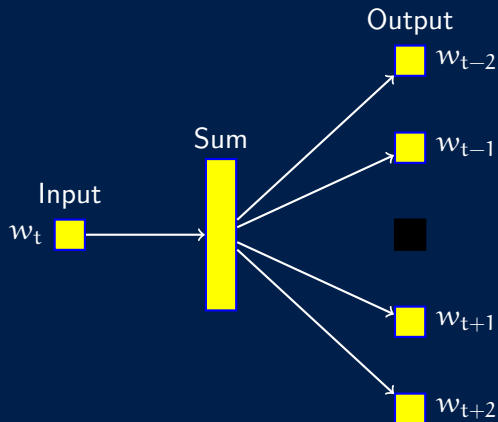


Figure: The CBOW architecture predicts the current word based on the context words of length n . Here the window size is 5

CBOW uses the sequence of words "Wish", "you", "a", "happy", "year" as a context and predicts or generates the central word "new"

- ▶ CBOW is used for learning the central word
- ▶ Maximize probability of word based on the word co-occurrences within a distance of n

SKIP GRAM MODEL



SG uses the central word "new" and predicts the context words "Wish", "you", "a", "happy", "year"

- ▶ SG is used to learn the context words given the central word
- ▶ Maximize probability of word based on the word co-occurrences within a distance of $[-n, +n]$ from the center word

Figure: The SG architecture predicts the one context word at a time based on the center word. Here the window size is 5

SOURCE PREPARATION FOR TRAINING

Source Text

Wish you more happy returns of the day→

Wish you more happy returns of the day→

Wish you more happy returns of the day→

Wish you more happy returns of the day→

Wish you more happy returns of the day→

Wish you many more happy returns of the day→

Wish you many more happy returns of the day →

Training Samples

(wish,you)

(wish,many)

(you,Wish)

(you,more),(you,happy)

(many,Wish),(many,you)

(many,more),(many,happy)

(more,many),(more,you)

(more,happy),(more,returns)

(happy,many),(happy,more)

(happy,returns),(happy,of)

(returns,more),(returns,happy)

(returns,of),(returns,the)

(of,happy),(of,returns)

(of,the),(of,day)

ONE-WORD LEARNING

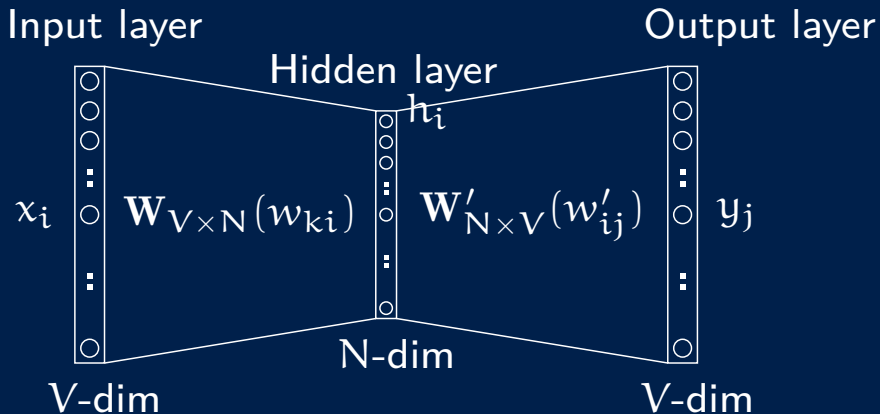


Figure: A CBOW model with only one word as input[1]. The layers are fully connected

Here the W and W' are learned

$$t^{\text{aback}} = \begin{pmatrix} 0 \\ 1 \\ \dots \\ 0 \\ \dots \\ 0 \\ 0 \end{pmatrix} \dots t^{\text{zoom}} = \begin{pmatrix} 0 \\ 0 \\ \dots \\ 0 \\ \dots \\ 1 \\ 0 \end{pmatrix} \quad t^{\text{zucchini}} = \begin{pmatrix} 0 \\ 0 \\ \dots \\ 0 \\ \dots \\ 0 \\ 1 \end{pmatrix}$$

$x_k = 1 \text{ and } x'_k = 0, \forall k' \neq k$

HIDDEN LAYER

This neural network is fully connected. Input to the network is a one-hot vector. \mathbf{v}_w^T is the N-dimensional vector representation of the word presented as input [2] [3].

$$\mathbf{h} = \mathbf{W}^T \mathbf{X} \quad (1)$$

Now \mathbf{v}_{w_I} of the matrix (W) is the vector representation of the input one-hot vector w_I . From (1), \mathbf{h} is a linear combination of input and weights. In the same way, we get a score for u_j

$$u_j = \mathbf{v}_{w_j}'^T \mathbf{h} = \mathbf{v}_{w_j}'^T \mathbf{v}_{w_I} \quad (2)$$

where \mathbf{v}_{w_I} is the vector representation of the input word w_I and \mathbf{v}_{w_j}' is the j^{th} column of (W')

OUTPUT LAYER

At the output layer, we apply the softmax to get the posterior distribution of the word(s). It is obtained by,

$$p(w_j | w_I) = y_j \quad (3)$$

where y_j is the output of the j^{th} unit in the output layer

$$y_j = \frac{\exp(u_j)}{\sum_{j'=1}^V \exp(u_{j'})} \quad (4)$$

$$= \frac{\exp(\mathbf{v}_{w_j}'^T \mathbf{v}_{w_I})}{\sum_{j'=1}^V \exp((\mathbf{v}_{w_{j'}}')^T \mathbf{v}_{w_I})} \quad (5)$$

where \mathbf{v}_w , \mathbf{v}'_w are the input vector (word vector) and output vector (feature vector) representations, of w_j and $w_{j'}$, respectively

UPDATE WEIGHTS - HIDDEN-OUTPUT LAYERS

The learning/training (through backpropagation [4]) objective is to maximize (5) or minimize the error between the target and the computed value of the target which is $y_j^* - t$ and t is same as the input vector, in this case. We use cross-entropy as it provides us with a good measure of "error distance"

$$\max p(w_o | w_I) = \max (\log(y_{j^*})) \text{ --- Maximize} \quad (6)$$

$$\text{---} E = u_j - \log(y_{j^*}) \text{ --- minimize} \quad (7)$$

$$= u_{j^*} - \log \sum_{j'=1}^V \exp(u_{j'}) \quad (8)$$

where

w_o is the output word and E is the loss function.

It is the special case of cross-entropy measurement between two probabilistic distributions u_{j^*} and $u_{j'}$

- ▶ $\log p(x)$ is well scaled
- ▶ Selection of step size is easier
- ▶ With $p(x)$ multiplication may yield to near zero causing *underflow*
- ▶ For better optimization, $\log p(x)$ is considered (multiplication \rightarrow addition)

UPDATE WEIGHTS (HO) - MINIMIZATION OF E

To minimize E, take the partial derivative of E with respect to j^{th} unit of u_j

$$\frac{\partial E}{\partial u_j} = y_j - t_j = e_j \quad (9)$$

where e_j is the prediction error. Taking partial derivative with respect to the hidden-output weights, we get,

$$\frac{\partial E}{\partial w'_{ij}} = \frac{\partial E}{\partial u_j} \cdot \frac{\partial u_j}{\partial w'_{ij}} = e_j \cdot h_i \quad (10)$$

Using the above equation (10),

$$w'_{ij}{}^{\text{new}} = w'_{ij}{}^{\text{old}} - \eta e_j \cdot h_i \text{ or} \quad (11)$$

$$v_{w_j}^{(\text{new})} = v_{w_j}^{(\text{old})} - \eta e_j \cdot \mathbf{h} \quad \text{for } j = 1, 2, 3, \dots, V \quad (12)$$

UPDATE INPUT TO HIDDEN WEIGHTS

Taking the derivative with respect to h_i , we get

$$\frac{\partial E}{\partial h_i} = \sum_{j=1}^V \frac{\partial E}{\partial u_j} \cdot \frac{\partial u_j}{\partial h_i} = \sum_{j=1}^V e_j \cdot w'_{ij} = \mathbf{E}\mathbf{H}_i \quad (13)$$

Taking the derivative with respect to w_{ki} , we get

$$\frac{\partial E}{\partial w_{ki}} = \frac{\partial E}{\partial h_i} \cdot \frac{\partial h_i}{\partial w_{ki}} = \mathbf{E}\mathbf{H}_i \cdot x_k \quad (14)$$

Now the weights are updated using

$$v_{wi}^{(\text{new})} = v_{wi}^{(\text{old})} - \eta \mathbf{E}\mathbf{H}^T \quad (15)$$

SOME INSIGHTS ON OUTPUT-HIDDEN-INPUT LAYER WEIGHT UPDATES

- ▶ The prediction error E propagates the weighted sum of all words in the vocabulary to every output vector v'_j
- ▶ The change in the input vector is defined by the output vector which in turn is updated due to the prediction error
- ▶ The model parameters accumulate the changes until the system reaches a state of equilibrium
- ▶ The rows in the Input-Hidden layer (v_j) stores the features of the words in the vocabulary V

MATRIX OPERATIONS

x	1	0	0	0	0	0	0	0	0	0
W	0.32	0.72	0.31	0.55	0.38	0.18	0.96	0.02	0.55	0.05
	0.25	0.90	0.61	0.01	0.23	0.91	0.75	0.71	0.22	0.56
	0.16	0.31	0.61	0.51	0.71	0.96	0.31	0.24	0.90	0.63
	0.66	0.88	0.15	0.47	0.68	0.76	0.37	0.69	0.40	0.94
	0.29	0.60	0.59	0.93	0.41	0.35	0.19	0.70	0.87	0.72

$h_{in} = W^T x$		0.32		0.72		0.31		0.05		0.32
		0.25		0.90		0.61		0.56		0.25
	$1 \times$	0.16	$+0 \times$	0.31	$+0 \times$	0.61	...	0.63	$=$	0.16
		0.66		0.88		0.15		0.94		0.66
		0.29		0.60		0.59		0.72		0.29

MATRIX OPERATIONS

W'										$\sigma(hin)$
0.67	0.32	0.46	0.61	0.02	0.85	0.04	0.69	0.58	0.65	0.58
0.56	0.96	0.36	0.81	0.62	0.49	0.99	0.15	0.41	0.35	0.56
0.34	0.27	0.32	0.84	0.29	0.51	0.24	0.56	0.42	0.59	0.54
0.67	0.15	0.98	0.26	0.23	0.91	0.92	0.54	0.86	0.06	0.66
0.21	0.77	0.39	0.05	0.61	0.69	0.52	0.83	0.17	0.68	0.57

$y^T = W'(W^T x)$	0.09	0.09	0.10	0.09	0.06	0.17	0.11	0.11	0.09	0.08
t	0	0	0	1	0	0	0	0	0	0
e	0.09	0.09	0.10	-0.91	0.06	0.17	0.11	0.11	0.09	0.08

cross entropy loss = 0.33

- ▶ Every element of the $y_{j^*} > 0$ and only one element of $t_j = 1$ for $j = j^*$.
 $e = y_j - t_j$
- ▶ If $e > 0$, then it is overestimated it
- ▶ $e < 0$ only when when $t_j = 1$. in such case, it is underestimated
- ▶ if $e \approx 0$, the the learning is complete and word vectors are learned. Now, the input vector is closer to the target vector.
 - ▶ In the case of CBOW model, the context vectors are now similar to the target vector
 - ▶ In the case of skipgram model, the input word found its similar words

MATRIX OPERATIONS

$$\text{matmul}(e, h)$$

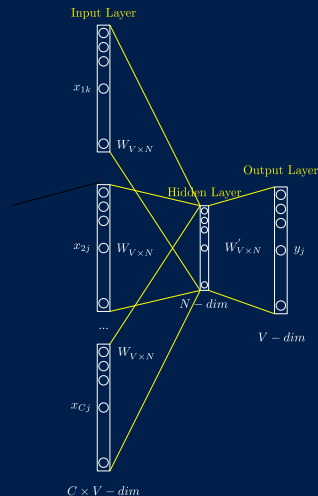
-0.53	0.05	0.06	0.05	0.03	0.10	0.06	0.06	0.05	0.05
-0.51	0.05	0.06	0.05	0.03	0.10	0.06	0.06	0.05	0.04
-0.49	0.05	0.05	0.05	0.03	0.09	0.06	0.06	0.05	0.04
-0.60	0.06	0.07	0.06	0.04	0.11	0.07	0.07	0.06	0.05
-0.52	0.05	0.06	0.05	0.03	0.10	0.06	0.06	0.05	0.05

$$W'^{\text{new}} = W'^{\text{old}} - \eta eh$$

0.667486	0.320777	0.461016	0.608380	0.023106	0.852884	0.041133	0.692782	0.583753	0.649416
0.667486	0.320777	0.461016	0.608380	0.023106	0.852884	0.041133	0.692782	0.583753	0.649416
0.555091	0.956016	0.364246	0.811703	0.615640	0.488893	0.992000	0.150837	0.408978	0.348874
0.341566	0.266564	0.322560	0.843457	0.292134	0.510061	0.236577	0.559690	0.424626	0.594844
0.671457	0.145615	0.982863	0.259544	0.227199	0.907842	0.920170	0.542782	0.858919	0.055786
0.207933	0.767069	0.393739	0.050281	0.614529	0.687904	0.519688	0.834847	0.169773	0.679568

CBOW MODEL FOR MULTIPLE WORDS

- ▶ C is the number of context words
- ▶ V is the size of the vocabulary
- ▶ h_i receives average of the vectors of the input context words
- ▶ Output vector v'_{wj} is the column vector in the \mathbf{W}' representing relationship between the context words and the target word
- ▶ Softmax is used for the output layer probability distribution for the target word



The CBOW Model

INPUT-HIDDEN WEIGHT VECTORS AND LOSS FUNCTION

The hidden units receive values from the linear combination of the context vectors and the weights

$$u_j = \mathbf{v}_{w_j}'^T \mathbf{h} = \mathbf{v}_{w_j}'^T \mathbf{v}_{w_I} \quad (16)$$

$$\begin{aligned} \mathbf{h} &= \frac{1}{C} \mathbf{W}^T (x_1 + x_2 + x_3 + \dots + x_C) \\ &= \frac{1}{C} (\mathbf{v}_{w1} + \mathbf{v}_{w2} + \mathbf{v}_{w3} + \dots + \mathbf{v}_{wC}) \end{aligned} \quad (17)$$

The equation for v_j' can be borrowed from (2) and E is

$$\begin{aligned} E &= -\log p(w_O | w_{I,1}, w_{I,2}, w_{I,3}, \dots, w_{I,C}) \\ &= -\mathbf{v}_{w_O}' \cdot \mathbf{h} + \log \sum_{j'=1}^V \exp \left(\mathbf{v}_{w_{j'}}'^T \cdot \mathbf{h} \right) \end{aligned} \quad (18)$$

UPDATE INPUT AND OUTPUT VECTORS

There is no change in the hidden-output weights² (12) as the computations remain the same. The new $\mathbf{v}_{\mathbf{w}_{I,c}}^{(\text{new})}$ is written as

$$\mathbf{v}_{\mathbf{w}_{I,c}}^{(\text{new})} = \mathbf{v}_{\mathbf{w}_{I,c}}^{(\text{old})} - \frac{1}{C} \cdot \eta \mathbf{E} \mathbf{H}^T, \text{ for } j = 1, 2, 3, \dots, C \quad (19)$$

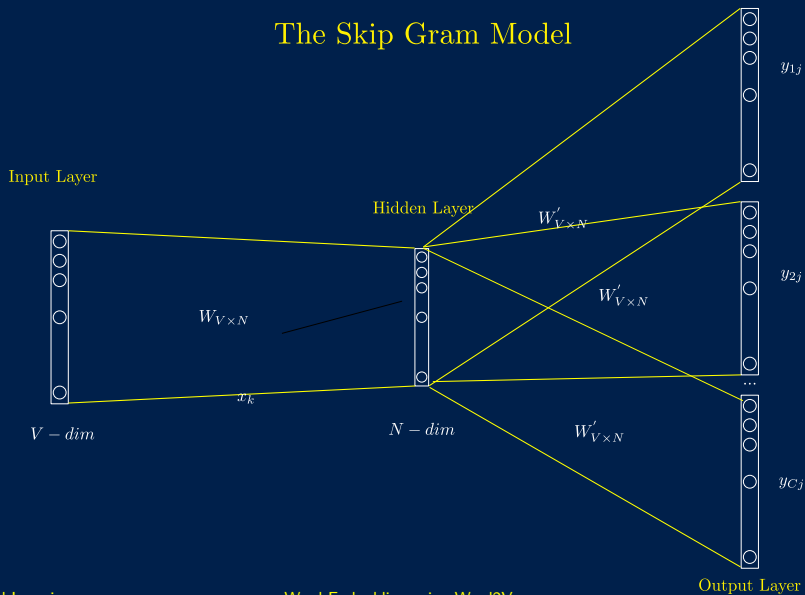
where η is the learning rate.

$$^2 \mathbf{v}_{\mathbf{w}_j}^{(\text{new})} = \mathbf{v}_{\mathbf{w}_j}^{(\text{old})} - \eta e_j \cdot \mathbf{h} \quad \text{for } j = 1, 2, 3, \dots, V$$

WHAT DOES IT LEARN?

- ▶ Distributed representation of words as vectors
- ▶ The learned vectors explicitly encode many linguistic regularities and patterns
- ▶ The learning should produce a similar word vectors for those words that appeared in similar context. How do we find out?
- ▶ Comparing the word vectors for similarity - Cosine similarity
- ▶ Has the learned word vectors address stemming? run, running, ran as similar?
 - ▶ He runs half-marathon
 - ▶ He ran half-marathon
 - ▶ He is running half-marathon
- ▶ How about car, cars, automobile?
- ▶ How about awesome, fantastic, great?

The Skip Gram Model



SKIP-GRAM MODEL - FORWARD AND BACK PROPAGATION UPDATE OF WEIGHTS

$$\text{Hidden layer } \mathbf{h} = W^T \mathbf{x} \quad (20)$$

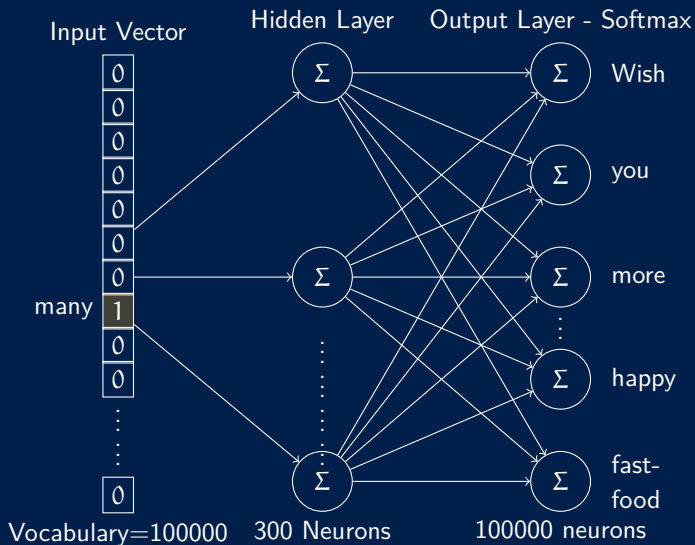
$$\text{Output neuron } u_{c,j} = \mathbf{v}'^T T_{w_j} \cdot \mathbf{h}, \text{ for } c = 1, 2, \dots, C \quad (21)$$

$$\text{cross entropy } E = - \sum_{c=1}^C u_j^* \cdot \log \sum_{j'=1}^{|V|} \exp(u_{j'}) \quad (22)$$

$$\text{Context weights } w_{ij}^{\text{new}} = w_{ij}^{\text{old}} - \eta \cdot E I_j \cdot h_i \quad (23)$$

$$\text{input weights } w_{ij}^{\text{new}} = w_{ij}^{\text{old}} - \eta E_j \cdot h_i \quad (24)$$

NEURAL NETWORK ARCHITECTURE - A SAMPLE



FORWARD PASS

```
def forward_pass(target_word_idx, W1, W2, vocab_size):  
    x = OHV(target_word_idx, vocab_size)  
    h = np.dot(x, W1)  
    u = np.dot(h, W2)  
    y_pred = softmax(u)  
    return x, h, y_pred
```

BACK PROPAGATION

```
def back_propagation(x, h, y_pred, context_word_idx,
                    input_weights, context_weights,
                    learning_rate):
    y_true = OHV(context_word_idx, vocab_size)
    e = y_pred - y_true
    dW2 = np.outer(h, e)
    dW1 = np.outer(x, np.dot(W2, e))
    input_weights -= learning_rate * dW1
```

TROUBLE WITH THE SIZE OF THE NETWORK

- ▶ All weights (output \rightarrow hidden) and (hidden \rightarrow input) are adjusted by taking a training sample so that the prediction cycle minimizes the loss function
- ▶ This amounts to updating all the weights in the neural network - amounts to several million weights for a network which has input neurons of size $|V| = 1\text{M}$, and hidden unit size as 300
- ▶ In addition, we should consider the several million training samples pairs

SOURCE PREPARATION FOR TRAINING

Source Text

Wish you many more happy returns of the day→

Wish you more happy returns of the day→

Wish you many more happy returns of the day→

Wish you many more happy returns of the day→

Wish you many more happy returns of the day→

Wish you many more happy returns of the day→

Wish you many more happy returns of the day →

Training Samples

(wish,you)

(wish,many)

(you,Wish)

(you,more),(you,happy)

(many,Wish),(many,you)

(many,more),(many,happy)

(more,many),(more,you)

(more,happy),(more,returns)

(happy,many),(happy,more)

(happy,returns),(happy,of)

(returns,more),(returns,happy)

(returns,of),(returns,the)

(of,happy),(of,returns)

(of,the),(of,day)

- ▶ The words (of, the) in the pairs (of, happy), (returns, the) do not give much information about the words happy and returns, respectively. Similarly, some pairs reappear with the order of the words switched.
- ▶ Some words could also be randomly removed from the based on the frequencies
- ▶ Words with less frequency or infrequent words appearing as context words could be discarded as they may not provide contextual information to the central word
- ▶ Subsampling is used to discard or downweight frequent words (like "the", "is", etc.) to focus training on less frequent but more informative words

SUB-SAMPLING IN WORD2VEC.C-GOOGLE

To counter the imbalance between the rare and frequent words, we used a simple subsampling approach: each word w_i in the training set is discarded with probability computed by the formula[5]

$$P(w_i) = 1 - \sqrt{\frac{1}{f(w_i)}} \quad (25)$$

Here is the code for sub-sampling used by `word2vec.c` that randomly removes a word from the sample using an algorithm similar to Linear Congruential generator (LCG) algorithm.

```
if (word == 0) break;
// The subsampling randomly discards frequent
//words while keeping the ranking same
if (sample > 0) {
    real ran = (sqrt(vocab[word].cn/(sample * train_words))+1)*
               (sample * train_words)/vocab[word].cn;
    next_random = next_random*(unsigned long long)25214903917+11;
    if (ran < (next_random & 0xFFFF) / (real)65536) continue;
}
```

SUB-SAMPLING IN WORD2VEC.C

- ▶ This calculates a random probability `ran` that is used to decide whether a word should be subsampled (discarded) or not.
- ▶ *sample*: Usually 10^{-5} . -making it very small makes the algorithm more aggressive in discarding frequent words.
- ▶ The term `vocab[word].cn/(sample*train_words)` measures the relative frequency of the word compared to the sampling threshold.
- ▶ The `sqrt()` function reduces the impact of very frequent words, allowing words that are just above the threshold to be less likely discarded.
- ▶ The entire expression generates a value `ran` that increases as the word becomes less frequent, meaning that less frequent words are more likely to be retained.

To deal with classification with multiple classes, softmax is very useful. If there are k classes in the data set, this activation function fits the classes in the range $[0,1]$ by calculating the probability. This is best suited for the finding the activation value of the neurons in the output layer. It is a normalized exponential function

$$P(C_k|x_j) = \frac{e^{a_j}}{\sum_k e^{a_k}}, \text{ where } k = 1, K \quad (26)$$

HIERARCHICAL SOFTMAX

- ▶ Has a flat hierarchy with a probability value for every output node of depth = 1
- ▶ Normalized over the probabilities of all $|V|$ words
- ▶ Error correction happens for every output \rightarrow hidden units
- ▶ Huge costs if the vocabulary size $|V|$ is of the order of several thousands
- ▶ Decompose the flat hierarchy into a binary tree
- ▶ Form a hierarchical description of a word as a sequence of $O(\log_2|V|)$ decisions and thereby reducing the computing complexity of Softmax - $O(|V|) \rightarrow O(\log_2(|V|))$
- ▶ Lay the words in a tree-based hierarchy - words as leaves
- ▶ Binary tree with $|V| - 1$ nodes for left (0) and right(1) traversal
- ▶ Every leaf represents the probability of the word

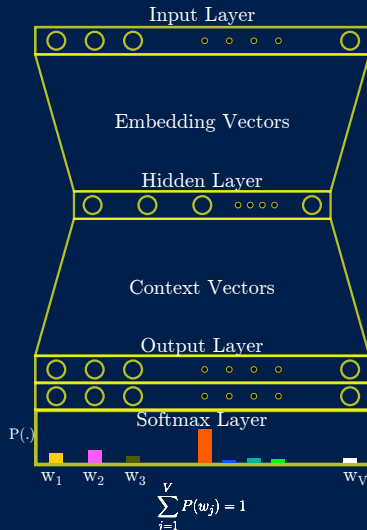
- ▶ Path length of a balanced Tree is $\log_2(|V|)$. If the $|V| = 1$ million words, then the path length = 19.9 bits/word
- ▶ Constructing an Huffman encoded-tree would help frequent words to have short unique binary codes
- ▶ Learn to take these probabilistic decisions instead of directly predicting each word's probability [6]
- ▶ Every intermediate node denotes the relative probabilities of its child nodes
- ▶ The path to reach every leaf (word) is unique
- ▶ H-Softmax in many cases increases the prediction speed by more than 50X times

SOFTMAX

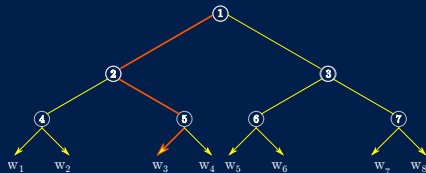
Softmax is a normalized exponential function

$$P(C_k|x_j) = \frac{e^{a_j}}{\sum_k e^{a_k}}, \text{ where } k=1, K \quad (27)$$

- ▶ Has a flat hierarchy with a probability value for every output node of depth = 1
- ▶ Normalized over the probabilities of all $|V|$ words
- ▶ Error correction happens for every output \rightarrow hidden units
- ▶ Huge costs if the vocabulary size $|V|$ is of the order of several thousands

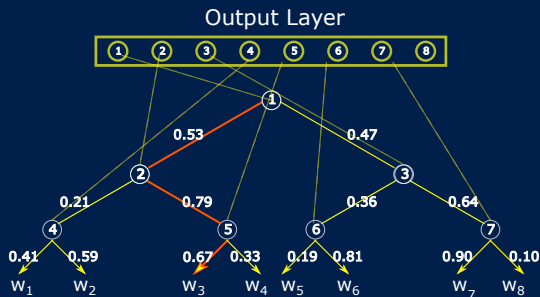


BALANCED BINARY TREE



- ▶ Move the words into a binary tree - depth depends on the vocabulary
 - ▶ The path to any word in the vocabulary is known
 - ▶ Traverse through the binary tree to reach any word
 - ▶ At every step, make a binary decision
- to reach the word
 - ▶ The length to reach any word in a balanced tree is $\log_2(|V|)$
 - ▶ Words could be arranged using
 - ▶ random order
 - ▶ IS-A relationship
 - ▶ TF-IDF frequency

BALANCED BINARY TREE WITH 8 WORDS



$P(w_i) = \prod_{j \in N_L} P(n(w, j))$, where N_L is the list of nodes to reach the word and

$$P(W) = \sum_{i=1}^V P(w_i) = 1$$

	w_1	w_2	w_3	w_4	w_5	w_6	w_7	w_8
$P(w_i)$	0.05	0.08	0.32	0.16	0.04	0.11	0.22	0.02

Thus the hierarchical Softmax is a well defined multinomial distribution among all words

HIERARCHICAL SOFTMAX - ADVANTAGES

- ▶ Decomposes the flat hierarchy into a binary tree
- ▶ The path to reach every leaf (word) is unique
- ▶ Lays the words in a tree-based hierarchy - words as leaves
- ▶ Binary tree with $|V| - 1$ nodes for left and right traversal
- ▶ Every intermediate node denotes the relative probabilities of its child nodes
- ▶ Every leaf represents the probability of the word

HIERARCHICAL SOFTMAX - ADVANTAGES

- ▶ Each node is indexed by a bit vector corresponding to the path from the root to the node
- ▶ Normalized values for the words are calculated without finding the probability for every word
- ▶ The entire vocabulary is partitioned into classes
- ▶ ANN learns to take these probabilistic decisions instead of directly predicting each word's probability³

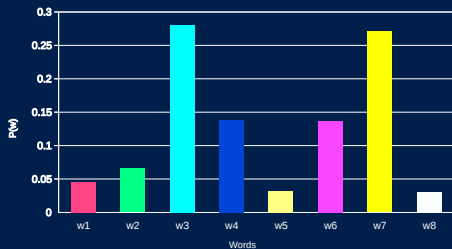
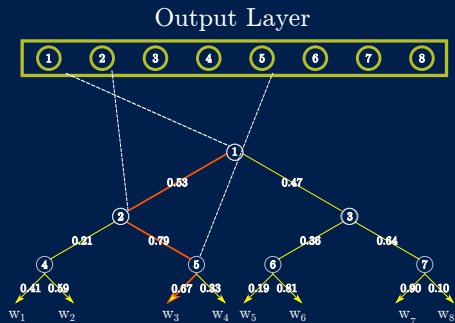
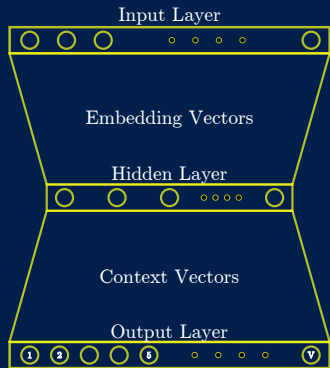
³Yoshua Bengio et al. "A Neural Probabilistic Model" In: J.Mach.Learn Res.3 (Mar.2003), pp 1137-1155. issn: 1532-4435

HIERARCHICAL SOFTMAX - ADVANTAGES

- ▶ Forms a hierarchical description of a word as a sequence of $O(\log_2|V|)$ decisions
- ▶ Reduces the computing complexity of Softmax - $O(|V|) \rightarrow O(\log_2(|V|))$
- ▶ Path length of a balanced Tree is $\log_2(|V|)$. If the $|V| = 1$ million words, then the path length = 19.9 bits/word
- ▶ A balanced binary tree should provide an exponential speed-up, on the order of $\frac{|V|}{\log_2(|V|)}$
- ▶ Constructing an Huffman encoded-tree would help frequent words to have short unique binary codes
- ▶ H-Softmax in many cases increases the prediction speed by more than 50X times

Note: Hierarchical softmax significantly reduces the computational complexity and it doesn't completely eliminate the need to compute probabilities for all words

WORD2VEC WITH HIERARCHICAL SOFTMAX



UPDATING WEIGHTS - 1/3

Let $L(w)$ be the number of nodes to traverse to the word from the root and $n(w, i)$ is the i^{th} node on this path and the associated vector in the context matrix is $v_{n(w, i)}$. $ch(n)$ is the child node [2][6][3][7]. Then the probability of word is

$$\begin{aligned} P(w|w_i) &= \prod_{j=1}^{L(w)-1} \sigma([n(w, j+1) = ch(n(w, j))] \cdot v_{n(w, j)}^T h) \\ &= \prod_{j=1}^{L(w)-1} \sigma([n(w, j+1) = ch(n(w, j))] \cdot v_{n(w, j)}^T v_{w_i}) \end{aligned} \quad (28)$$

$$\text{where } [x] = \begin{cases} 1, & \text{if } x \text{ is true} \\ -1, & \text{otherwise} \end{cases} \quad \text{and } \sigma(\cdot) \text{ is the sigmoid function}$$

if the child node $ch(n(w, j))$ is left of the parent node, then the term $[n(w, j+1) = ch(n(w, j))]$ is 1, and equal to -1, if the path goes to the right.

Since the sum of the probabilities of at the node is 1, we can prove that

$$\sigma(\mathbf{v}_n^T \mathbf{v}_{w_i}) + \sigma(-\mathbf{v}_n^T \mathbf{v}_{w_i}) = 1, \text{ since } 1 - \sigma(x) = \sigma(-x) \quad (29)$$

Example

$$P(w|w_i) = \sigma(\mathbf{v}_{n(w,j)}^T \mathbf{v}_{w_i}) \cdot \sigma(-\mathbf{v}_{n(w,j)}^T \mathbf{v}_{w_i}) \cdot \sigma(\mathbf{v}_{n(w,j)}^T \mathbf{v}_{w_i})$$

To train the model, we need to minimize the negative log likelihood $-\log P(w|w_i)$

$$E = - \sum_{j=1}^{L(w)-1} \log \sigma([\cdot] u'_j), \text{ where } u_j = \mathbf{v}'_j \cdot \mathbf{h} \quad (30)$$

$$(31)$$

where $t_j = 1$, if $[n(w, j+1) = \text{ch}(n(w, j))] = 1$ and $t_j = 0$ otherwise

$$\frac{\partial E}{\partial u_j} = \sigma(u_j - 1)[.] \quad (32)$$

$$= \begin{cases} \sigma(u_j) - 1 & \text{for } [.] = 1 \\ \sigma(u_j) & \text{for } [.] = -1 \end{cases} \quad (33)$$

$$= \sigma(u_j) - t_j \quad (34)$$

where $t_j = 1$, if

$[n(w, j + 1) = \text{ch}(n(w, j))] = 1$, else $t_j = 0$

$$\frac{\partial E}{\partial v'_j} = \sigma(v'_j \cdot h) - t_j \quad (35)$$

$$v'_j{}^{\text{new}} = v'_j{}^{\text{old}} - \eta(\sigma(v'_j \cdot h) - t_j) \cdot h \quad (36)$$

$$\frac{\partial E}{\partial h} = \sum_{j=1}^{L(W)-1} \frac{\partial E}{\partial v'_j \cdot h} \cdot \frac{\partial v'_j \cdot h}{\partial h} = EH \quad (37)$$

$$v_{w_i}^{\text{new}} = v_{w_i}^{\text{old}} - \eta \cdot EH^T \quad (38)$$

NEGATIVE SAMPLING

- ▶ The size of the network is proportional to the size of the vocabulary V . For every training cycle of input, the every weight in the network needs to be updated
- ▶ For every training cycle, Softmax function computes the sum of the output neuron values
- ▶ Cost of updating all the weights in the fully connected network is very high
- ▶ Is it possible to change only a small percentage of the weights?
- ▶ Select a small number of *negative* words
- ▶ While updating the weights, these samples output zero while the positive sample(s) will retain its value
- ▶ During the backpropagation, the weights related to the negative and positive words are changed and the rest will remain untouched for the current update
- ▶ This reduces drastically the computation 😊

SELECTING A NEGATIVE SAMPLE

$$P(w_i) = \frac{f(w_i)^{\frac{3}{4}}}{\sum_{j=0}^n f(w_j)^{\frac{3}{4}}} \quad (39)$$

Using the frequency of words directly for negative sampling would bias the model towards frequent words. Negative sampling does favor more frequent words, but it uses a smoothed frequency distribution with the $\frac{3}{4}$ power to ensure a good balance between frequent and rare words.

NEGATIVE SAMPLING - AN EXAMPLE

- ▶ Modify a small percentage of the weights, rather than all of them
- ▶ Consider **chalk** and **blackboard** pair for one-word learning one-word learning.
 - ▶ We want blackboard to have one ($= 1$) at the output layer
 - ▶ Traditionally, we want rest of the words to have zero ($= 0$)
 - ▶ with Negative sampling, we pick a random set of words. Let us say we pick up words coffee, car, elephant, tower as negative samples. Instead of making every word zero, we turn the negative values to be equal to zero and rest unchanged

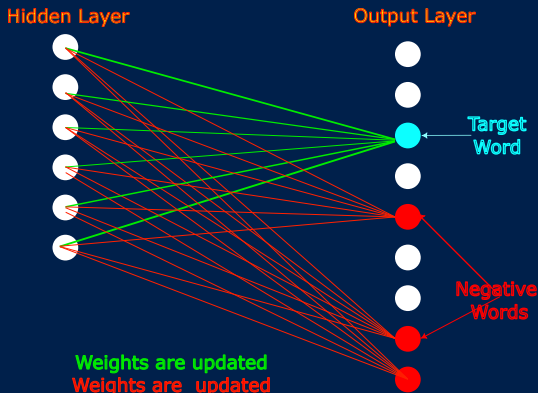
Logistic regression is used loss function used in when Negative Sampling is used:

$$E = - \left(\log(\sigma(z_{\text{pos}})) + \sum_{i=1}^k \log(1 - \sigma(z_{\text{neg}_i})) \right) \quad (40)$$

where z_{pos} is the input for the positive word, aiming for a value close to 1 and z_{neg_i} (for $i = 1$ to k) represents the inputs for each of the k negative words, aiming for values close to 0.

ADVANTAGES OF NEGATIVE SAMPLING

- ▶ **Faster training time** - especially for large datasets. It allows scaling of the
- ▶ **Better accuracy** Improves accuracy by avoiding over fitting
- ▶ **Optimized utilization of memory**



RESULTS - SMILAR WORDS FOR Virus

Vocab size	Words in the corpus
637722	222502540

Word	Similarity
virus,	0.889620
viral	0.785719
(herpesvirus)	0.764385
avirus	0.759567
fluav)	0.757418
polio-virus	0.724740
⋮	⋮
(vsv;	0.723436
(denv-2)	0.722825
(cowpox)	0.717185
⋮	⋮

ANALOGIES - virus : covid :: ransomware : ?

Word	Distance
scams	0.605626
cybercrime	0.592376
denial-of-service	0.590397
⋮	⋮
privacy	0.589996
cyberattack	0.586179
frauds	0.585077
wartime	0.584866
(ddos)	0.582940
⋮	⋮

BUILDING WORD2VEC FROM SCRATCH

1. Data Preparation

- ▶ Preprocess the text corpus
- ▶ Create a vocabulary of unique words.
- ▶ Construct a reverse lookup to map words to their indices

2. Generate Training Data

- ▶ Initialize a sliding window of size w
- ▶ For each word in the corpus:
 - ▶ Extract the center word and context words
 - ▶ Add the pair to the training data

3. Implement Skip-Gram Model

- ▶ Define the following layers:
 - ▶ *Input layer*: One-hot encoded center word
 - ▶ *Hidden layer*: Word embeddings
 - ▶ *Softmax layer*: predict context words

4. Training

- ▶ Initialize the model parameters
- ▶ For each epoch:
 - ▶ For every training data, compute the loss
 - ▶ Update the model parameters

5. Negative Sampling

- ▶ Use negative sampling speed up the training
 - ▶ Optimize the model to distinguish context words from negative words

6. Training Loop

- ▶ Implement a training loop to iterate through the following steps:
 - ▶ Compute the gradients of the model parameters
 - ▶ Update the model parameters based on the gradients

7. Save Word Embeddings

- ▶ Save the word embeddings to a file for future use

8. Evaluation

- ▶ Evaluate the word embeddings using tasks like word similarity, word analogy, or text classification.

9. Optimization and Performance

- ▶ Optimize your code for performance, as training Word2Vec can be computationally intensive
 - ▶ Use multi-threading or distributed computing if needed

COMPLEXITIES OF THESE MODELS

Task	Complexity
Skip-gram (Hierarchical Softmax)	$O(T \cdot N \cdot C \cdot \log V)$
Skip-gram (Negative Sampling)	$O(T \cdot N \cdot C \cdot k)$
CBOW (Hierarchical Softmax)	$O(T \cdot N \cdot C \cdot \log V)$
CBOW (Negative Sampling)	$O(T \cdot N \cdot C \cdot k)$
Cosine Similarity	$O(N)$
Analogy Task	$O(V \cdot N)$

where

T is total number of words in the corpus

V is vocabulary size

N is embedding size

C is context window size

k is number of negative samples

LIMITATIONS

- ▶ Separate training is required for phrases
- ▶ Embeddings are learned based on a small local window surrounding words - good and bad share the almost the same embedding
- ▶ Does not address polysemy
- ▶ Does not use frequencies of term co-occurrences

SOURCE CODE FOR WORD2VEC

GitHub link to W2V

REFERENCES I

- [1] Xin Rong. *word2vec Parameter Learning Explained*. 2014. arXiv: [1411.2738](https://arxiv.org/abs/1411.2738). URL: https://sites.socsci.uci.edu/~lpearl/courses/readings/Rong2014_Word2Vec.pdf.
- [2] Jeffrey A. Dean Tomas Mikolov Kai Chen Gregory S. Corrado. U.S. pat. US9037464B1. May 2015.
- [3] Tomas Mikolov et al. “Distributed Representations of Words and Phrases and Their Compositionality”. In: *Proceedings of the 26th International Conference on Neural Information Processing Systems - Volume 2*. NIPS’13. Lake Tahoe, Nevada: Curran Associates Inc., 2013, pp. 3111–3119. URL: <http://dl.acm.org/citation.cfm?id=2999792.2999959>.
- [4] David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. “Learning representations by back-propagating errors”. In: *nature* 323.6088 (1986), pp. 533–536.

- [5] Quoc Le and Tomas Mikolov. “Distributed representations of sentences and documents”. In: *International conference on machine learning*. 2014, pp. 1188–1196. URL: <http://proceedings.mlr.press/v32/le14.pdf>.
- [6] Yoshua Bengio et al. “A Neural Probabilistic Language Model”. In: *Journal of Machine Learning Research* 3 (Mar. 2003), pp. 1137–1155. ISSN: 1532-4435. URL: <http://dl.acm.org/citation.cfm?id=944919.944966>.
- [7] Andriy Mnih and Geoffrey Hinton. “A Scalable Hierarchical Distributed Language Model”. In: *Proceedings of the 21st International Conference on Neural Information Processing Systems*. NIPS’08. Vancouver, British Columbia, Canada: Curran Associates Inc., 2008, pp. 1081–1088. ISBN: 978-1-6056-0-949-2. URL: <http://dl.acm.org/citation.cfm?id=2981780.2981915>.