

**f** (<https://www.facebook.com/AnalyticsVidhya>)

**t** (<https://twitter.com/analyticsvidhya>)

**g+** (<https://plus.google.com/+Analyticsvidhya/posts>)

**in** (<https://www.linkedin.com/groups/Analytics-Vidhya-Learn-everything-about-5057165>)



(<https://www.analyticsvidhya.com/datahacksummit/>)

Home (<https://www.analyticsvidhya.com/>) > Deep Learning (<https://www.analyticsvidhya.com/blog/category/deep-learning/>) >

# An Intuitive Understanding of Word Embeddings: From Count Vectors to Word2Vec

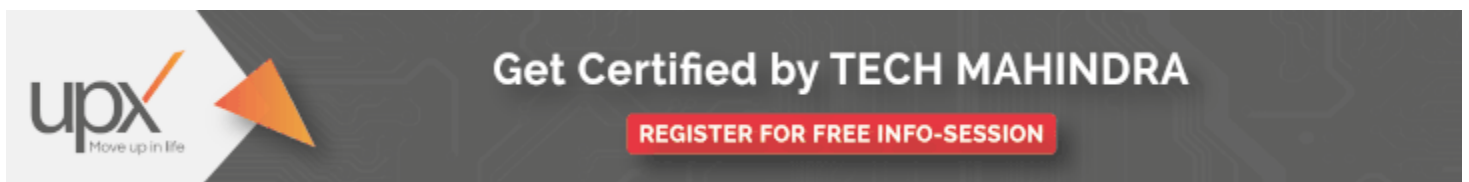
DEEP LEARNING (<https://www.analyticsvidhya.com/blog/category/deep-learning/>) MACHINE LEARNING

(<https://www.analyticsvidhya.com/blog/category/machine-learning/>) NLP

(<https://www.analyticsvidhya.com/blog/category/nlp/>) PYTHON (<https://www.analyticsvidhya.com/blog/category/python-2/>)

2/)

sharer.php?u=<https://www.analyticsvidhya.com/blog/2017/06/word-embeddings-count-understanding%20of%20%20Word%20Embeddings:%20From%20Count%20Vectors%20to%20Word2Vec>) **t** ([https://plus.google.com/share?url=https://www.analyticsvidhya.com/blog/2017/06/word-embeddings-count-word2veec/](https://twitter.com/home?lang%20of%20%20Word%20Embeddings:%20From%20Count%20Vectors%20to%20Word2Vec+https://www.analyticsvidhya.com/blog/2017/06/word-tps://plus.google.com/share?url=https://www.analyticsvidhya.com/blog/2017/06/word-embeddings-count-word2veec/)) **p**  
1/?url=<https://www.analyticsvidhya.com/blog/2017/06/word-embeddings-count-word2veec/&media=https://s3-ap-south-1.amazonaws.com/av-blog-062705/Word-re%20Understanding%20of%20%20Word%20Embeddings:%20From%20Count%20Vectors%20to%20Word2Vec>)

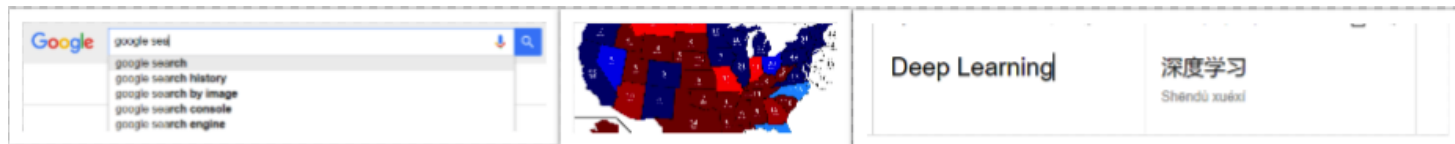


([http://events.upxacademy.com/infosession?utm\\_source=MLWeek-AVBnr&utm\\_medium=Banner&utm\\_campaign=Info-Session](http://events.upxacademy.com/infosession?utm_source=MLWeek-AVBnr&utm_medium=Banner&utm_campaign=Info-Session))

# Introduction

Before we start, have a look at the below examples.

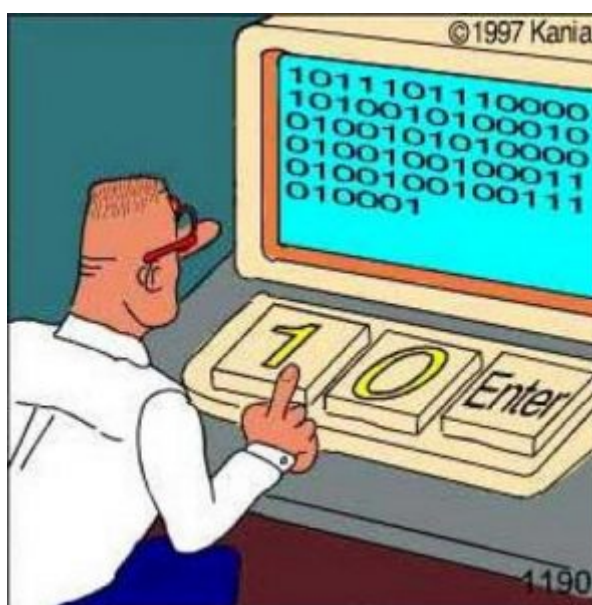
1. You open Google and search for a news article on the ongoing Champions trophy and get hundreds of search results in return about it.
2. Nate silver analysed millions of tweets and correctly predicted the results of 49 out of 50 states in 2008 U.S Presidential Elections.
3. You type a sentence in google translate in English and get an Equivalent Chinese conversion.



So what do the above examples have in common?

You possibly guessed it right – **TEXT processing**. All the above three scenarios deal with humongous amount of text to perform different range of tasks like clustering in the google search example, Machine translation in the second and classification in the third.

Humans can deal with text format quite intuitively but provided we have millions of documents being generated in a single day, we cannot have humans performing the above the three tasks. It is neither scalable nor effective.



So, how do we make computers of today perform clustering, classification etc on a text data since we know that they are generally inefficient at handling and processing strings or texts for any fruitful outputs?

Sure, a computer can match two strings and tell you whether they are same or not. But how do we make computers tell you about football or Ronaldo when you search for Messi? How do you make a computer understand that "Apple" in "Apple is a tasty fruit" is a fruit that can be eaten and not a company?

The answer to the above questions lie in creating a representation for words that capture their *meanings, semantic relationships* and the different types of contexts they are used in.

And all of these are implemented by using Word Embeddings or numerical representations of texts so that computers may handle them.

Below, we will see formally what are Word Embeddings and their different types and how we can actually implement them to perform the tasks like returning efficient Google search results.

## Table of Contents

1. What are Word Embeddings?
2. Different types of Word Embeddings
  - 2.1 Frequency based Embedding
    - 2.1.1 Count Vectors
    - 2.1.2 TF-IDF
    - 2.1.3 Co-Occurrence Matrix
  - 2.2 Prediction based Embedding
    - 2.2.1 CBOW
    - 2.2.2 Skip-Gram
3. Word Embeddings use case scenarios(what all can be done using word embeddings? eg: similarity, odd one out etc.)
4. Using pre-trained Word Vectors
5. Training your own Word Vectors
6. End Notes

## 1. What are Word Embeddings?

In very simplistic terms, Word Embeddings are the texts converted into numbers and there may be different numerical representations of the same text. But before we dive into the details of Word Embeddings, the following question should be asked – Why do we need Word Embeddings?

As it turns out, many Machine Learning algorithms and almost all Deep Learning Architectures are incapable of processing *strings* or *plain text* in their raw form. They require numbers as inputs to perform any sort of job, be it classification, regression etc. in broad terms. And with the huge amount of data that is present in the text format, it is imperative to extract knowledge out of it and build applications. Some real world applications of text applications are – sentiment analysis of reviews by Amazon etc., document or news classification or clustering by Google etc.

Let us now define Word Embeddings formally. A Word Embedding format generally tries to map a word using a dictionary to a vector. Let us break this sentence down into finer details to have a clear view.

Take a look at this example – **sentence**= " Word Embeddings are Word converted into numbers "

A *word* in this **sentence** may be "Embeddings" or "numbers " etc.

A *dictionary* may be the list of all unique words in the **sentence**. So, a dictionary may look like – ['Word','Embeddings','are','Converted','into','numbers']

A *vector* representation of a word may be a one-hot encoded vector where 1 stands for the position where the word exists and 0 everywhere else. The vector representation of "numbers" in this format according to the above dictionary is [0,0,0,0,0,1] and of converted is [0,0,0,1,0,0].

This is just a very simple method to represent a word in the vector form. Let us look at different types of Word Embeddings or Word Vectors and their advantages and disadvantages over the rest.

## 2. Different types of Word Embeddings

The different types of word embeddings can be broadly classified into two categories-

1. Frequency based Embedding
2. Prediction based Embedding

Let us try to understand each of these methods in detail.

## 2.1 Frequency based Embedding

There are generally three types of vectors that we encounter under this category.

1. Count Vector
2. TF-IDF Vector
3. Co-Occurrence Vector

Let us look into each of these vectorization methods in detail.

### 2.1.1 Count Vector

Consider a Corpus C of D documents  $\{d_1, d_2, \dots, d_D\}$  and N unique tokens extracted out of the corpus C. The N tokens will form our dictionary and the size of the Count Vector matrix M will be given by  $D \times N$ . Each row in the matrix M contains the frequency of tokens in document D(i).

Let us understand this using a simple example.

D1: He is a lazy boy. She is also lazy.

D2: Neeraj is a lazy person.

The dictionary created may be a list of unique tokens(words) in the corpus =  
['He','She','lazy','boy','Neeraj','person']

Here,  $D=2$ ,  $N=6$

The count matrix M of size  $2 \times 6$  will be represented as –

	He	She	lazy	boy	Neeraj	person
D1	1	1	2	1	0	0
D2	0	0	1	0	1	1

Now, a column can also be understood as word vector for the corresponding word in the matrix M. For example, the word vector for 'lazy' in the above matrix is [2,1] and so on. Here, the *rows* correspond to the *documents* in the corpus and the *columns* correspond to the *tokens* in the dictionary. The second row in the above matrix may be read as – D2 contains 'lazy': once, 'Neeraj': once and 'person' once.

Now there may be quite a few variations while preparing the above matrix M. The variations will be generally in-

1. The way dictionary is prepared.

Why? Because in real world applications we might have a corpus which contains millions of documents. And with millions of document, we can extract hundreds of millions of unique words. So basically, the matrix that will be prepared like above will be a very sparse one and inefficient for any computation. So an alternative to using every unique word as a dictionary element would be to pick say top 10,000 words based on frequency and then prepare a dictionary.

2. The way count is taken for each word.

We may either take the frequency (number of times a word has appeared in the document) or the presence(has the word appeared in the document?) to be the entry in the count matrix M. But generally, frequency method is preferred over the latter.

Below is a representational image of the matrix M for easy understanding.

	Document 1	Document 2	Document 3	Document 4	Document 5	Document 6	Document 7	Document 8
Term(s) 1	10	0	1	0	0	0	0	2
Term(s) 2	0	2	0	0	0	18	0	2
Term(s) 3	0	0	0	0	0	0	0	2
Term(s) 4	6	0	0	4	6	0	0	0
Term(s) 5	0	0	0	0	0	0	0	2
Term(s) 6	0	0	1	0	0	1	0	0
Term(s) 7	0	1	8	0	0	0	0	0
Term(s) 8	0	0	0	0	0	3	0	0

← Word Vector (Passage Vector)

Document Vector

## 2.1.2 TF-IDF vectorization

This is another method which is based on the frequency method but it is different to the count vectorization in the sense that it takes into account not just the occurrence of a word in a single document but in the entire corpus. So, what is the rationale behind this? Let us try to understand.

Common words like 'is', 'the', 'a' etc. tend to appear quite frequently in comparison to the words which are important to a document. For example, a document **A** on Lionel Messi is going to contain more occurrences of the word "Messi" in comparison to other documents. But common words like "the" etc. are also going to be present in higher frequency in almost every document.

Ideally, what we would want is to down weight the common words occurring in almost all documents and give more importance to words that appear in a subset of documents.

TF-IDF works by penalising these common words by assigning them lower weights while giving importance to words like Messi in a particular document.

So, how exactly does TF-IDF work?

Consider the below sample table which gives the count of terms(tokens/words) in two documents.

Document 1		Document 2	
Term	Count	Term	Count
This	1	This	1
is	1	is	2
about	2	about	1
Messi	4	Tf-idf	1

Now, let us define a few terms related to TF-IDF.

TF = (Number of times term t appears in a document)/(Number of terms in the document)

So, TF(This,Document1) =  $1/8$

TF(This, Document2)= $1/5$

It denotes the contribution of the word to the document i.e words relevant to the document should be frequent. eg: A document about Messi should contain the word 'Messi' in large number.

IDF =  $\log(N/n)$ , where, N is the number of documents and n is the number of documents a term t has appeared in.

where N is the number of documents and n is the number of documents a term t has appeared in.

So, IDF(This) =  $\log(2/2) = 0$ .

So, how do we explain the reasoning behind IDF? Ideally, if a word has appeared in all the document, then probably that word is not relevant to a particular document. But if it has appeared in a subset of documents then probably the word is of some relevance to the documents it is present in.

Let us compute IDF for the word 'Messi'.

$$\text{IDF}(\text{Messi}) = \log(2/1) = 0.301.$$

Now, let us compare the TF-IDF for a common word 'This' and a word 'Messi' which seems to be of relevance to Document 1.

$$\text{TF-IDF}(\text{This}, \text{Document1}) = (1/8) * (0) = 0$$

$$\text{TF-IDF}(\text{This}, \text{Document2}) = (1/5) * (0) = 0$$

$$\text{TF-IDF}(\text{Messi}, \text{Document1}) = (4/8) * 0.301 = 0.15$$

As, you can see for Document1, TF-IDF method heavily penalises the word 'This' but assigns greater weight to 'Messi'. So, this may be understood as 'Messi' is an important word for Document1 from the context of the entire corpus.

### 2.1.3 Co-Occurrence Matrix with a fixed context window

**The big idea** – Similar words tend to occur together and will have similar context for example – Apple is a fruit. Mango is a fruit.

Apple and mango tend to have a similar context i.e fruit.

Before I dive into the details of how a co-occurrence matrix is constructed, there are two concepts that need to be clarified – Co-Occurrence and Context Window.

**Co-occurrence** – For a given corpus, the co-occurrence of a pair of words say  $w_1$  and  $w_2$  is the number of times they have appeared together in a Context Window.

**Context Window** – Context window is specified by a number and the direction. So what does a context window of 2 (around) means? Let us see an example below,



Quick	Brown	Fox	Jump	Over	The	Lazy	Dog
-------	-------	-----	------	------	-----	------	-----

The green words are a 2 (around) context window for the word 'Fox' and for calculating the co-occurrence only these words will be counted. Let us see context window for the word 'Over'.

Quick	Brown	Fox	Jump	Over	The	Lazy	Dog
-------	-------	-----	------	------	-----	------	-----

Now, let us take an example corpus to calculate a co-occurrence matrix.

Corpus = He is not lazy. He is intelligent. He is smart.

	He	is	not	lazy	intelligent	smart
He	0	4	2	1	2	1
is	4	0	1	2	2	1
not	2	1	0	1	0	0
lazy	1	2	1	0	0	0
intelligent	2	2	0	0	0	0
smart	1	1	0	0	0	0

Let us understand this co-occurrence matrix by seeing two examples in the table above. Red and the blue box.

Red box- It is the number of times 'He' and 'is' have appeared in the context window 2 and it can be seen that the count turns out to be 4. The below table will help you visualise the count.

He	is	not	lazy	He	is	intelligent	He	is	smart
He	is	not	lazy	He	is	intelligent	He	is	smart
He	is	not	lazy	He	is	intelligent	He	is	smart
He	is	not	lazy	He	is	intelligent	He	is	smart

while the word 'lazy' has never appeared with 'intelligent' in the context window and therefore has been assigned 0 in the blue box.

## Variations of Co-occurrence Matrix

Let's say there are  $V$  unique words in the corpus. So Vocabulary size =  $V$ . The columns of the Co-occurrence matrix form the *context words*. The different variations of Co-Occurrence Matrix are-

1. A co-occurrence matrix of size  $V \times V$ . Now, for even a decent corpus  $V$  gets very large and difficult to handle. So generally, this architecture is never preferred in practice.
2. A co-occurrence matrix of size  $V \times N$  where  $N$  is a subset of  $V$  and can be obtained by removing irrelevant words like stopwords etc. for example. This is still very large and presents computational difficulties.

But, remember this co-occurrence matrix is not the word vector representation that is generally used. Instead, this Co-occurrence matrix is decomposed using techniques like PCA, SVD etc. into factors and combination of these factors forms the word vector representation.

Let me illustrate this more clearly. For example, you perform PCA on the above matrix of size  $V \times V$ . You will obtain  $V$  principal components. You can choose  $k$  components out of these  $V$  components. So, the new matrix will be of the form  $V \times k$ .

And, a single word, instead of being represented in  $V$  dimensions will be represented in  $k$  dimensions while still capturing almost the same semantic meaning.  $k$  is generally of the order of hundreds.

So, what PCA does at the back is decompose Co-Occurrence matrix into three matrices,  $U$ ,  $S$  and  $V$  where  $U$  and  $V$  are both orthogonal matrices. What is of importance is that dot product of  $U$  and  $S$  gives the word vector representation and  $V$  gives the word context representation.

$$\begin{pmatrix} \hat{X} \\ \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & \\ \vdots & \vdots & \ddots & \\ x_{m1} & & & x_{mn} \end{pmatrix} \\ m \times n \end{pmatrix} \approx \begin{pmatrix} U \\ \begin{pmatrix} u_{11} & \dots & u_{1r} \\ \vdots & \ddots & \\ u_{m1} & & u_{mr} \end{pmatrix} \\ m \times r \end{pmatrix} \begin{pmatrix} S \\ \begin{pmatrix} s_{11} & 0 & \dots \\ 0 & \ddots & \\ \vdots & & s_{rr} \end{pmatrix} \\ r \times r \end{pmatrix} \begin{pmatrix} V^T \\ \begin{pmatrix} v_{11} & \dots & v_{1n} \\ \vdots & \ddots & \\ v_{r1} & & v_{rn} \end{pmatrix} \\ r \times n \end{pmatrix}$$

## Advantages of Co-occurrence Matrix

1. It preserves the semantic relationship between words. i.e man and woman tend to be closer than man and apple.
2. It uses SVD at its core, which produces more accurate word vector representations than existing methods.
3. It uses factorization which is a well-defined problem and can be efficiently solved.
4. It has to be computed once and can be used anytime once computed. In this sense, it is faster in comparison to others.

### Disadvantages of Co-Occurrence Matrix

1. It requires huge memory to store the co-occurrence matrix.  
But, this problem can be circumvented by factorizing the matrix out of the system for example in Hadoop clusters etc. and can be saved.

## 2.2 Prediction based Vector

**Pre-requisite:** This section assumes that you have a working knowledge of how a neural network works and the mechanisms by which weights in an NN are updated. If you are new to Neural Network, I would suggest you go through [this awesome article](https://www.analyticsvidhya.com/blog/2017/05/neural-network-from-scratch-in-python-and-r/) (https://www.analyticsvidhya.com/blog/2017/05/neural-network-from-scratch-in-python-and-r/) by Sunil to gain a very good understanding of how NN works.

So far, we have seen deterministic methods to determine word vectors. But these methods proved to be limited in their word representations until Mitolov etc. el introduced word2vec to the NLP community. These methods were prediction based in the sense that they provided probabilities to the words and proved to be state of the art for tasks like word analogies and word similarities. They were also able to achieve tasks like King -man +woman = Queen, which was considered a result almost magical. So let us look at the word2vec model used as of today to generate word vectors.

Word2vec is not a single algorithm but a combination of two techniques – CBOW(Continuous bag of words) and Skip-gram model. Both of these are shallow neural networks which map word(s) to the target variable which is also a word(s). Both of these techniques learn weights which act as word vector representations. Let us discuss both these methods separately and gain intuition into their working.

### 2.2.1 CBOW (Continuous Bag of words)

The way CBOW work is that it tends to predict the probability of a word given a context. A context may be a single word or a group of words. But for simplicity, I will take a single context word and try to predict a single target word.

Suppose, we have a corpus C = "Hey, this is sample corpus using only one context word." and we have defined a context window of 1. This corpus may be converted into a training set for a CBOW model as follow. The input is shown below. The matrix on the right in the below image contains the one-hot encoded from of the input on the left.

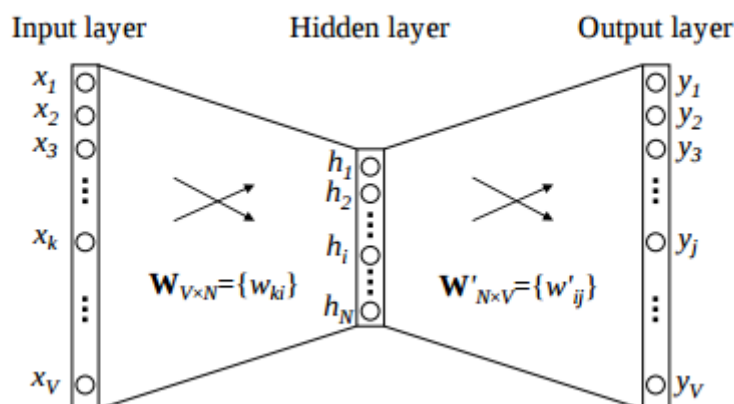
Input	Output		Hey	This	is	sample	corpus	using	only	one	context	word
Hey	this	Datapoint 1	1	0	0	0	0	0	0	0	0	0
this	hey	Datapoint 2	0	1	0	0	0	0	0	0	0	0
is	this	Datapoint 3	0	0	1	0	0	0	0	0	0	0
is	sample	Datapoint 4	0	0	1	0	0	0	0	0	0	0
sample	is	Datapoint 5	0	0	0	1	0	0	0	0	0	0
sample	corpus	Datapoint 6	0	0	0	1	0	0	0	0	0	0
corpus	sample	Datapoint 7	0	0	0	0	1	0	0	0	0	0
corpus	using	Datapoint 8	0	0	0	0	1	0	0	0	0	0
using	corpus	Datapoint 9	0	0	0	0	0	1	0	0	0	0
using	only	Datapoint 10	0	0	0	0	0	1	0	0	0	0
only	using	Datapoint 11	0	0	0	0	0	0	1	0	0	0
only	one	Datapoint 12	0	0	0	0	0	0	1	0	0	0
one	only	Datapoint 13	0	0	0	0	0	0	0	1	0	0
one	context	Datapoint 14	0	0	0	0	0	0	0	1	0	0
context	one	Datapoint 15	0	0	0	0	0	0	0	0	1	0
context	word	Datapoint 16	0	0	0	0	0	0	0	0	1	0
word	context	Datapoint 17	0	0	0	0	0	0	0	0	0	1

The target for a single datapoint say Datapoint 4 is shown as below

Hey	this	is	sample	corpus	using	only	one	context	word
0	0	0	1	0	0	0	0	0	0

This matrix shown in the above image is sent into a shallow neural network with three layers: an input layer, a hidden layer and an output layer. The output layer is a softmax layer which is used to sum the probabilities obtained in the output layer to 1. Now let us see how the forward propagation will work to calculate the hidden layer activation.

Let us first see a diagrammatic representation of the CBOW model.



The matrix representation of the above image for a single data point is below.

Context										Input-Hidden Weight				Hidden Activation			
C1	this	0	1	0	0	0	0	0	0	0	0	0	0	5	6	7	8
										1	2	3	4				
										5	6	7	8				
										9	10	11	12				
										13	14	15	16				
										17	18	19	20				
										21	22	23	24				
										25	26	27	28				
										29	30	31	32				
										33	34	35	36				
										37	38	39	40				

The flow is as follows:

1. The input layer and the target, both are one-hot encoded of size  $[1 \times V]$ . Here  $V=10$  in the above example.
2. There are two sets of weights. one is between the input and the hidden layer and second between hidden and output layer.  
Input-Hidden layer matrix size  $= [V \times N]$ , hidden-Output layer matrix size  $= [N \times V]$ : Where  $N$  is the number of dimensions we choose to represent our word in. It is arbitrary and a hyper-parameter for a Neural Network. Also,  $N$  is the number of neurons in the hidden layer. Here,  $N=4$ .
3. There is a no activation function between any layers. (and by no, i mean linear activation)
4. The input is multiplied by the input-hidden weights and called hidden activation. It is simply the corresponding row in the input-hidden matrix copied.
5. The hidden input gets multiplied by hidden-output weights and output is calculated.
6. Error between output and target is calculated and propagated back to re-adjust the weights.
7. The weight between the hidden layer and the output layer is taken as the word vector representation of the word.

We saw the above steps for a single context word. Now, what about if we have multiple context words? The image below describes the architecture for multiple context words.



The image above takes 3 context words and predicts the probability of a target word. The input can be assumed as taking three one-hot encoded vectors in the input layer as shown above in red, blue and green.

So, the input layer will have 3  $[1 \times V]$  Vectors in the input as shown above and 1  $[1 \times V]$  in the output layer. Rest of the architecture is same as for a 1-context CBOW.

The steps remain the same, only the calculation of hidden activation changes. Instead of just copying the corresponding rows of the input-hidden weight matrix to the hidden layer, an average is taken over all the corresponding rows of the matrix. We can understand this with the above figure. The average vector calculated becomes the hidden activation. So, if we have three context words for a single target word, we will have three initial hidden activations which are then averaged element-wise to obtain the final activation.

In both a single context word and multiple context word, I have shown the images till the calculation of the hidden activations since this is the part where CBOW differs from a simple MLP network. The steps after the calculation of hidden layer are same as that of the MLP as mentioned

in this article – Understanding and Coding Neural Networks from scratch

(<https://www.analyticsvidhya.com/blog/2017/05/neural-network-from-scratch-in-python-and-r/>).

The differences between MLP and CBOW are mentioned below for clarification:

1. The objective function in MLP is a MSE(mean square error) whereas in CBOW it is negative log likelihood of a word given a set of context i.e  $-\log(p(w_o|w_i))$ , where  $p(w_o|w_i)$  is given as

$$p(w_o|w_i) = \frac{\exp(v'_{w_o} \top v_{w_i})}{\sum_{w=1}^W \exp(v'_w \top v_{w_i})}$$

$w_o$  : output word

$w_i$ : context words

2. The gradient of error with respect to hidden-output weights and input-hidden weights are different since MLP has sigmoid activations(generally) but CBOW has linear activations. The method however to calculate the gradient is same as an MLP.

### Advantages of CBOW:

1. Being probabilistic in nature, it is supposed to perform superior to deterministic methods(generally).
2. It is low on memory. It does not need to have huge RAM requirements like that of co-occurrence matrix where it needs to store three huge matrices.

### Disadvantages of CBOW:

1. CBOW takes the average of the context of a word (as seen above in calculation of hidden activation). For example, Apple can be both a fruit and a company but CBOW takes an average of both the contexts and places it in between a cluster for fruits and companies.
2. Training a CBOW from scratch can take forever if not properly optimized.

## 2.2.2 Skip – Gram model

Skip – gram follows the same topology as of CBOW. It just flips CBOW’s architecture on its head. The aim of skip-gram is to predict the context given a word. Let us take the same corpus that we built our CBOW model on. C="Hey, this is sample corpus using only one context word." Let us construct the training data.

Input	Output(Context1)	Output(Context2)
Hey	this	<padding>
this	Hey	is
is	this	sample
sample	is	corpus
corpus	sample	corpus
using	corpus	only
only	using	one
one	only	context
context	one	word
word	context	<padding>

The input vector for skip-gram is going to be similar to a 1-context CBOW model. Also, the calculations up to hidden layer activations are going to be the same. The difference will be in the target variable. Since we have defined a context window of 1 on both the sides, there will be **"two" one hot encoded target variables** and **"two" corresponding outputs** as can be seen by the blue section in the image.

Two separate errors are calculated with respect to the two target variables and the two error vectors obtained are added element-wise to obtain a final error vector which is propagated back to update the weights.

The weights between the input and the hidden layer are taken as the word vector representation after training. The loss function or the objective is of the same type as of the CBOW model.

The skip-gram architecture is shown below.



[illegible]

Let us break down the above image.

Input layer size –  $[1 \times V]$ , Input hidden weight matrix size –  $[V \times N]$ , Number of neurons in hidden layer –  $N$ , Hidden-Output weight matrix size –  $[N \times V]$ , Output layer size –  $C [1 \times V]$

In the above example, C is the number of context words=2, V= 10, N=4

1. The row in red is the hidden activation corresponding to the input one-hot encoded vector. It is basically the corresponding row of input-hidden matrix copied.
2. The yellow matrix is the weight between the hidden layer and the output layer.
3. The blue matrix is obtained by the matrix multiplication of hidden activation and the hidden output weights. There will be two rows calculated for two target(context) words.
4. Each row of the blue matrix is converted into its softmax probabilities individually as shown in the green box.

5. The grey matrix contains the one hot encoded vectors of the two context words(target).
6. Error is calculated by subtracting the first row of the grey matrix(target) from the first row of the green matrix(output) element-wise. This is repeated for the next row. Therefore, for **n** target context words, we will have **n** error vectors.
7. Element-wise sum is taken over all the error vectors to obtain a final error vector.
8. This error vector is propagated back to update the weights.

## Advantages of Skip-Gram Model

1. Skip-gram model can capture two semantics for a single word. i.e it will have two vector representations of Apple. One for the company and other for the fruit.
2. Skip-gram with negative sub-sampling outperforms every other method generally.

This (<http://bit.ly/wevi-online>) is an excellent interactive tool to visualise CBOW and skip gram in action. I would suggest you to really go through this link for a better understanding.

## 3. Word Embeddings use case scenarios

Since word embeddings or word Vectors are numerical representations of contextual similarities between words, they can be manipulated and made to perform amazing tasks like-

1. Finding the degree of similarity between two words.  

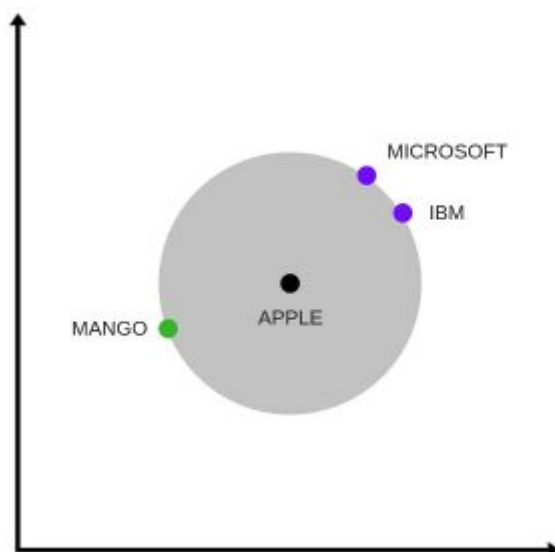
```
model.similarity('woman', 'man')
0.73723527
```
2. Finding odd one out.  

```
model.doesnt_match('breakfast cereal dinner lunch';.split())
'cereal'
```
3. Amazing things like woman+king-man =queen  

```
model.most_similar(positive=['woman', 'king'], negative=['man'], topn=1)
queen: 0.508
```
4. Probability of a text under the model  

```
model.score(['The fox jumped over the lazy dog'.split()])
0.21
```

Below is one interesting visualisation of word2vec.



The above image is a t-SNE representation of word vectors in 2 dimension and you can see that two contexts of apple have been captured. One is a fruit and the other company.

5. It can be used to perform Machine Translation.



The above graph is a bilingual embedding with chinese in green and english in yellow. If we know the words having similar meanings in chinese and english, the above bilingual embedding can be used to translate one language into the other.

## 4. Using pre-trained word vectors

We are going to use google's pre-trained model. It contains word vectors for a vocabulary of 3 million words trained on around 100 billion words from the google news dataset. The download link for the model is this

(<https://drive.google.com/file/d/0B7XkCwpl5KDYNlNUTTlSS21pQmM/edit>). Beware it is a 1.5 GB download.

```
from gensim.models import Word2Vec

#loading the downloaded model
model = Word2Vec.load_word2vec_format('GoogleNews-vectors-negative300.bin',
binary=True, norm_only=True)

#the model is loaded. It can be used to perform all of the tasks mentioned above.

# getting word vectors of a word
dog = model['dog']

#performing king queen magic
print(model.most_similar(positive=['woman', 'king'], negative=['man']))

#picking odd one out
print(model.doesnt_match("breakfast cereal dinner lunch".split()))

#printing similarity index
print(model.similarity('woman', 'man'))
```

## 5. Training your own word vectors

We will be training our own word2vec on a custom corpus. For training the model we will be using gensim and the steps are illustrated as below.

word2Vec requires that a format of list of list for training where every document is contained in a list and every list contains list of tokens of that documents. I won't be covering the pre-processing part here. So let's take an example list of list to train our word2vec model.

```
sentence=[['Neeraj','Boy'],['Sarwan','is'],['good','boy']]

#training word2vec on 3 sentences
model = gensim.models.Word2Vec(sentence, min_count = 1,size=300,workers=4 )
```

Let us try to understand the parameters of this model.

sentence – list of list of our corpus

min\_count=1 -the threshold value for the words. Word with frequency greater than this only are going to be included into the model.

size=300 – the number of dimensions in which we wish to represent our word. This is the size of the word vector.

workers=4 – used for parallelization

```
#using the model
```

```
#The new trained model can be used similar to the pre-trained ones.
```

```
#printing similarity index
```

```
print(model.similarity('woman', 'man'))
```

## 6. End Notes

Word Embeddings is an active research area trying to figure out better word representations than the existing ones. But, with time they have grown large in number and more complex. This article was aimed at simplifying some of the workings of these embedding models without carrying the mathematical overhead. If you feel think that I was able to clear some of your confusion, comment below. Any changes or suggestions would be welcomed.

---

### Share this:

 (<https://www.analyticsvidhya.com/blog/2017/06/word-embeddings-count-word2veec/?share=linkedin&nb=1>)

9

 (<https://www.analyticsvidhya.com/blog/2017/06/word-embeddings-count-word2veec/?share=facebook&nb=1>)

14

 (<https://www.analyticsvidhya.com/blog/2017/06/word-embeddings-count-word2veec/?share=google-plus-1&nb=1>)

 (<https://www.analyticsvidhya.com/blog/2017/06/word-embeddings-count-word2veec/?share=twitter&nb=1>)

 (<https://www.analyticsvidhya.com/blog/2017/06/word-embeddings-count-word2veec/?share=pocket&nb=1>)

 (<https://www.analyticsvidhya.com/blog/2017/06/word-embeddings-count-word2veec/?share=reddit&nb=1>)