

<b>WordNet</b> (Discrete representation)	<b>One-hot representation</b> (count vectorizing)	<b>Bag of words</b>	<b>TF-IDF</b>
+ Good as a resource	+ Simplest to implement	+ captures word frequency	+ weights to represent whether a word is significant
<ul style="list-style-type: none"> <li>- Missing <b>nuances</b> (e.g. synonyms)</li> <li>- Missing <b>new words</b></li> <li>- Needs <b>human</b> labor to create &amp; adapt</li> <li>- Hard to compute accurate <b>word similarity</b></li> <li>- Regards words as atomic symbols</li> </ul>	<ul style="list-style-type: none"> <li>- <b>Unordered</b>, context of words is lost</li> <li>- vector rep is in binary form, no <b>freq</b> info considered</li> <li>- Really <b>huge</b> and <b>sparse</b> vectors</li> </ul>	<ul style="list-style-type: none"> <li>- Can't capture word order</li> <li>- Can't capture similarities (e.g. boy, girl, giraffe. BOW cannot say who are more similar)</li> <li>- High dimensional &amp; sparse</li> </ul>	<ul style="list-style-type: none"> <li>- high dimensional</li> <li>- don't capture semantic relatedness</li> </ul>

<http://hunterheidenreich.com/blog/intro-to-word-embeddings/>

## Dense Vector Embeddings

(Singular value decomposition etc)

<b>Window-based Co-occurrence Matrix</b> Word-doc CM gives general topics - WBCM -> symmetric (irrelevant whether left or right context)	<b>t-SNE (For visualization) use after SVD</b> (rep multi-dim vectors in 2D space; suited for visualization) (optimizes st similar points in high-dim space are still closeby in low-dim space)
+ captures both syntactic (POS) & semantic info	
<ul style="list-style-type: none"> <li>- Increase size with vocab</li> <li>- Very high dim</li> <li>- Subsequent classification model will have sparsity issue (less robust)</li> </ul> <p>Sol: Low dimensional vectors          Store 'most' of the impt info in a fixed small # dim          Reduce dimensionality via <b>SVD</b></p>	<b>SVD</b> <ul style="list-style-type: none"> <li>- For <math>n \times m</math> matrix: cost = <math>O(mn^2)</math>; bad when dealing with millions of words / docs</li> <li>- Hard to incorporate new words/doc</li> </ul> <p>Let's directly learn low-dim word vectors!</p>

## Word representations

Traditional Method - Bag of Words Model	Word Embeddings
<ul style="list-style-type: none"> <li>• <b>one hot</b> encoding</li> <li>• Each word in the vocabulary is represented by one bit position in a HUGE vector.</li> <li>• For example, if we have a vocabulary of 10000 words, and "Hello" is the 4<sup>th</sup> word in the dictionary, it would be represented by: 0 0 0 1 0 0 . . . . . 0 0 0 0</li> <li>• Context information is not utilized</li> </ul>	<ul style="list-style-type: none"> <li>• Each word is a <b>point in <math>n</math>-dimensional space</b>, represented by a <b>vector of length <math>n</math></b>, (<math>n \sim 100 - 300</math>)</li> <li>• Trained by using big text dataset</li> <li>• For example, "Hello" might be represented as : [0.4, -0.11, 0.55, 0.3 . . . 0.1, 0.02]</li> <li>• Dimensions are basically projections along different axes, more of a mathematical concept.</li> </ul>

## Directly learning word vectors

<b>Word2vec</b> Efficient embedding model Instead of capturing co-occurrence counts directly, <b>predict surround words of each word</b> <ul style="list-style-type: none"><li>- Maximise <math>P(\text{context} \mid \text{focus word})</math></li><li>- 2 options (both 1 layer): CBOW, Skip-gram (SG more popular)</li><li>- <b>CBOW</b>: predict word given context, <b>SG</b> opposite</li></ul>	<b>Negative Sampling as objective</b> <ul style="list-style-type: none"><li>- new formulation of objective that implements binary logistic regression to classify between data &amp; noise samples</li></ul>	<b>GloVe</b> (efficient count-based model) Combines elements from 2 main word embedding models Instead of learning raw co-occurrence probs, <b>learn ratios of these co-occurrence prob</b>  <b>Noise points</b> : words that don't help us distinguish b/w i and j (use ratio b/w CO prob to filter out these noise points)  $P(k \mid j) / P(k \mid i)$ with different k words. <b>If value is high -&gt; k is useful to distinguish i and j</b>
To train this neural network, we need to <ol style="list-style-type: none"><li>1. Decide window size (dependent on dataset size / quality)</li><li>2. Decide embedding size (#nodes in network)</li><li>3. Size of vocab</li></ol> Training obj? Softmax cross-entropy? -> hard to compute Solution: <b>Negative Sampling</b>	Results: <ul style="list-style-type: none"><li>+ Fast &amp; Accurate predictive model</li><li>+ Captures semantic content</li><li>+ Pretrained models available</li><li>+ Super fast to train</li></ul> <i>(these models try to encode some meaning in the words so that you can use it in your eventual classification model)</i>	<ul style="list-style-type: none"><li>+ Fast training</li><li>+ Scalable to huge corpus, good perf</li><li>+ Easier to parallelize than word2vec, but comparable accuracy</li><li>+ Ratio of co-occurrence probabilities can encode meaning</li></ul>

### Negative Sampling for Word2Vec:

Instead of considering all words in vocabulary, consider a few 'negative samples' for each 'positive sample'. Randomly sample 5 negative (false) contexts for each positive (correct) (word, context) in the dataset.

### Evaluation of word vectors

1. Intrinsic evaluation: it encodes semantic information (e.g. similarity)
2. Extrinsic evaluation: is it useful for other NLP tasks

Word similarity task, Word analogy task

Extrinsic evaluation: Part of Speech tagging, Named Entity recognition

### Application for data science:

- Input to other models
- Sentiment analysis

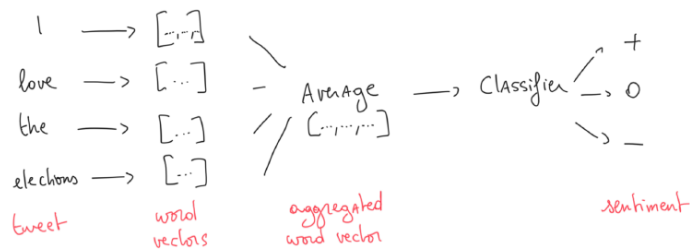
### Sentiment Analysis:

1. Preprocessing: remove URLs, weird characters, code tags, etc

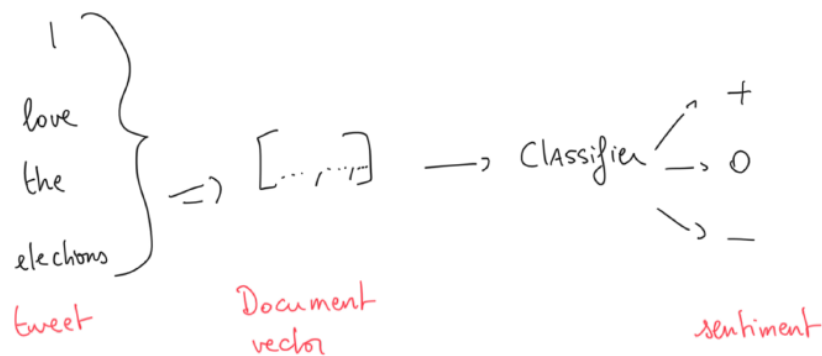
2. Tokenize (i.e. split text into tokens like words)
3. Get embeddings for words (average the vectors to overcome variable length input)

Does not work well...

Example : tweet sentiment prediction



→ problem: does not work super well..



→ Doc2vec

Use **Doc2vec**

- > Generalise word2vec to whole documents
- > Provides fixed-length vector, one vector for your document
- + Faster and consumes less memory than word2vec
- + Typically obtains better results