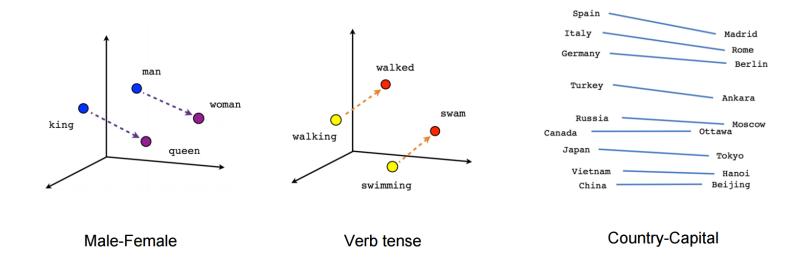
Vector Representation of Text

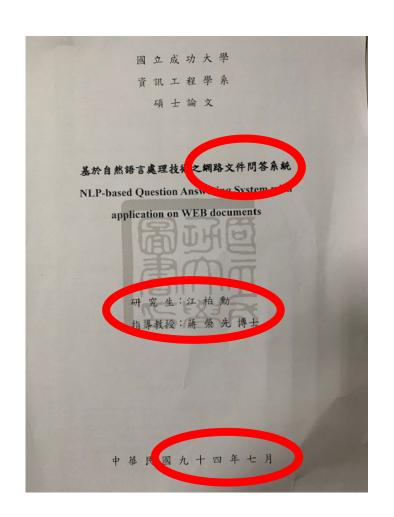
Word Embedding Technique (word2vec)

Word to vector (word2vector)

- The more often two words co-occur, the closer their vectors will be
- Two words have close meanings if their local neighborhoods are similar

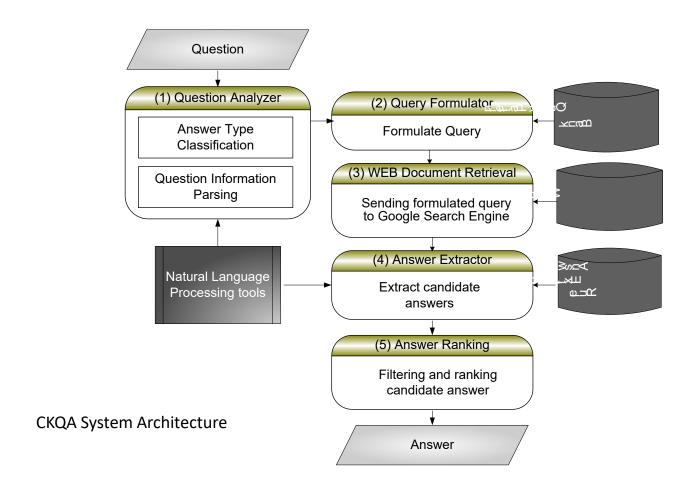


我們IIR Lab的團隊



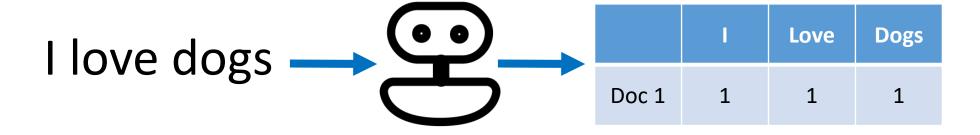


研究架構與做法

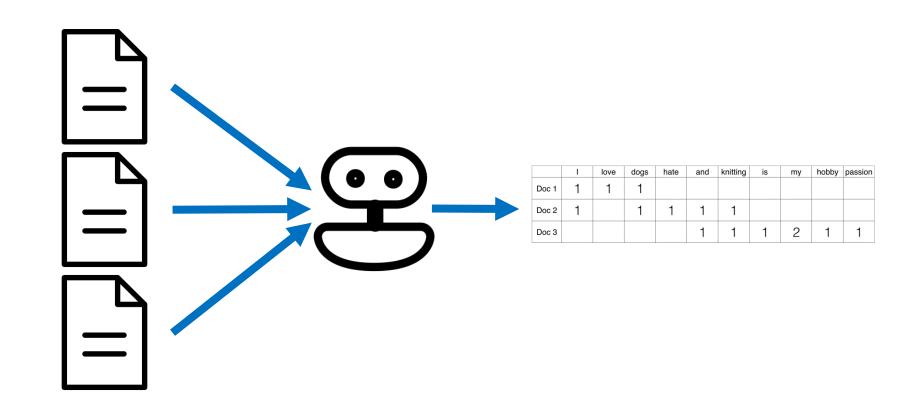




Bag of Words



Bag of Words

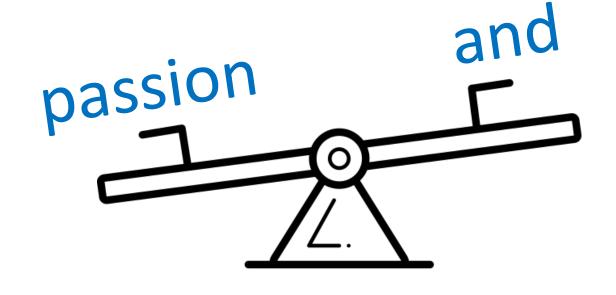


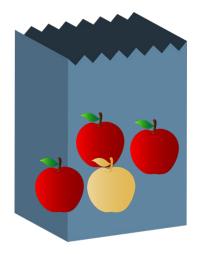
	l	love	dogs	hate	and	knitting	is	my	hobby	passion
Doc 1	1	1	1							
Doc 2	1		1	1	5	1				
Doc 3					1	1	1	2	1	1

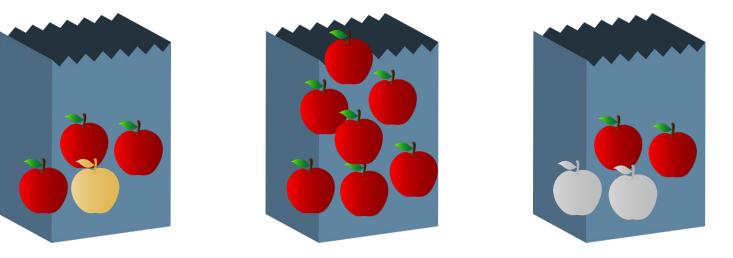


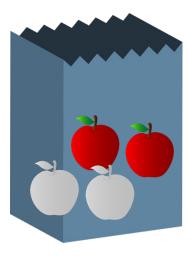


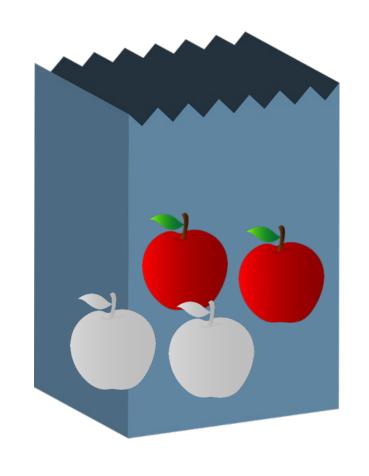




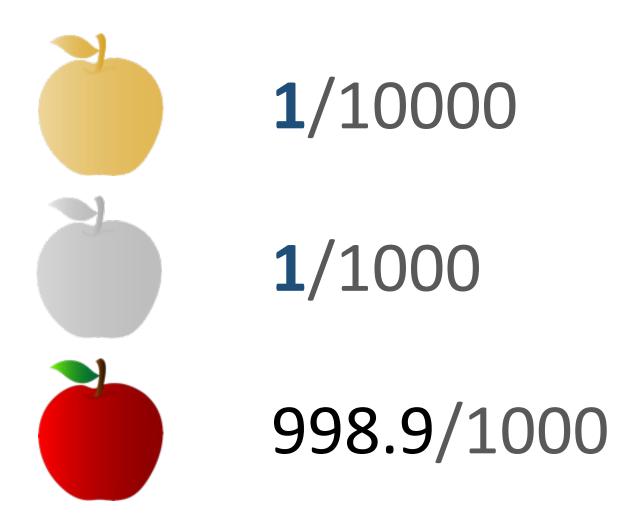


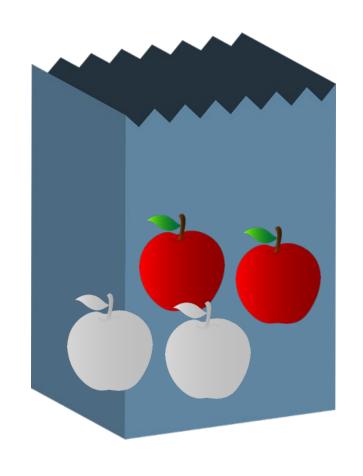




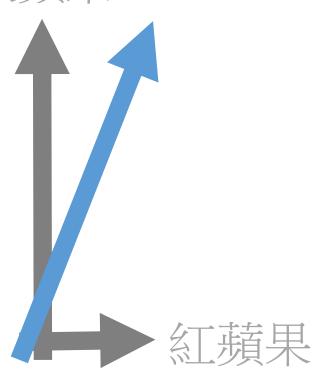


銀蘋果 對於這袋子 有多重要?





銀蘋果



$$w_{x,y} = tf_{x,y} \times log(\frac{N}{df_x})$$

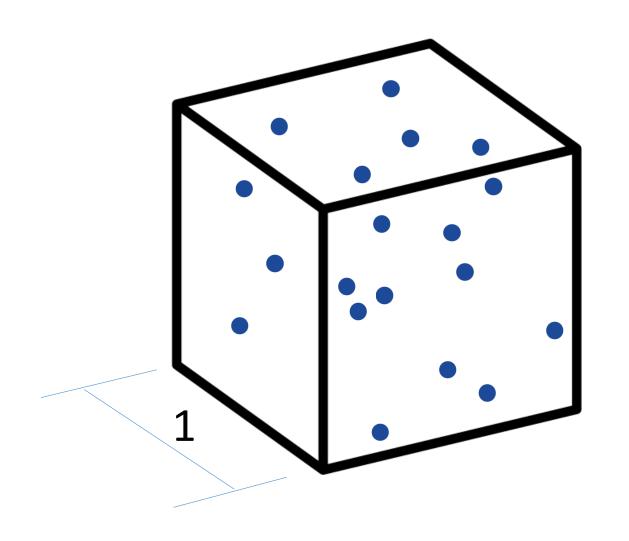
TF-IDF

Term x within document y

 $tf_{x,y}$ = frequency of x in y df_x = number of documents containing x N = total number of documents

Word embedding

keras.layers.Embedding

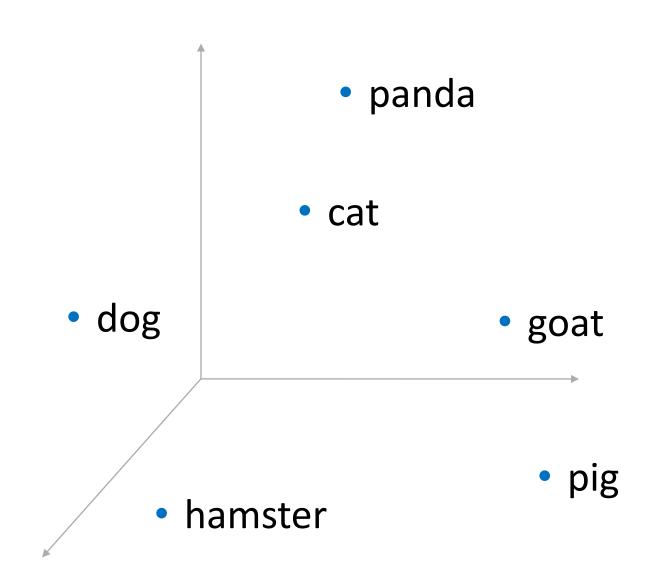


passion

x: 0.119

y: 0.212

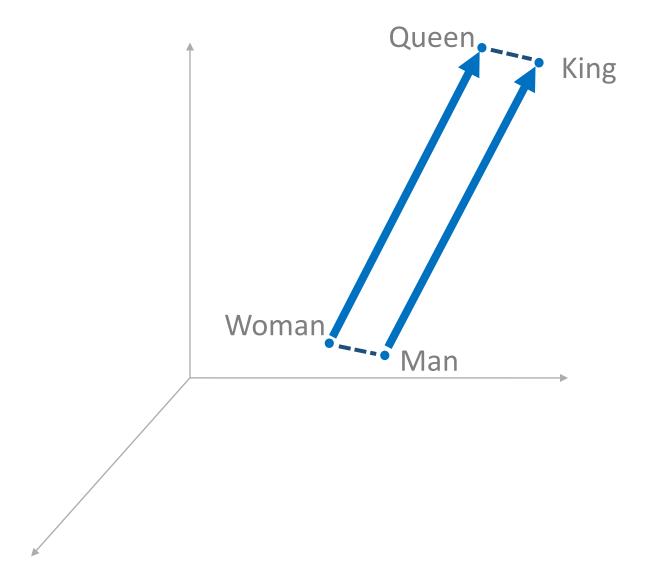
z: 0.010



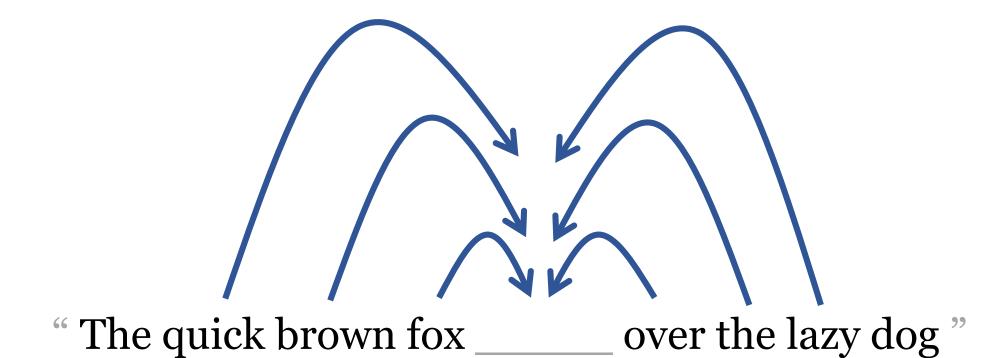
Word2Vec

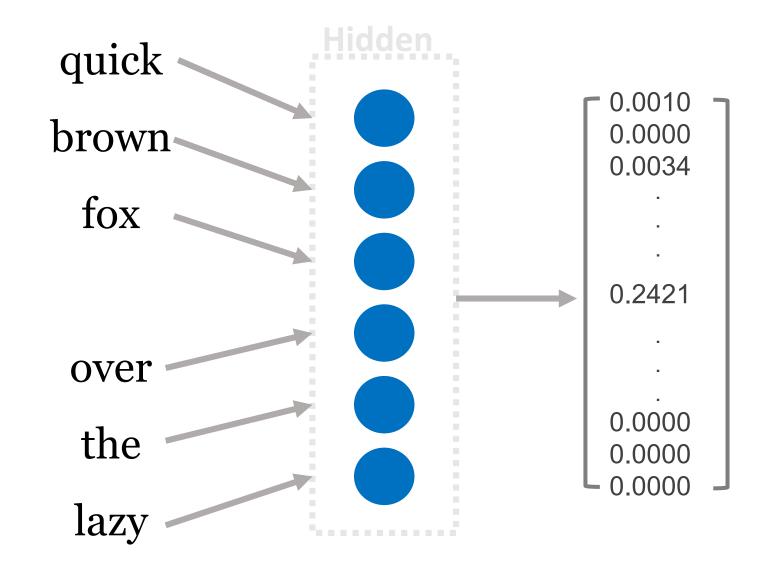
```
panda
cat
goat
dog
pig
hamster
```

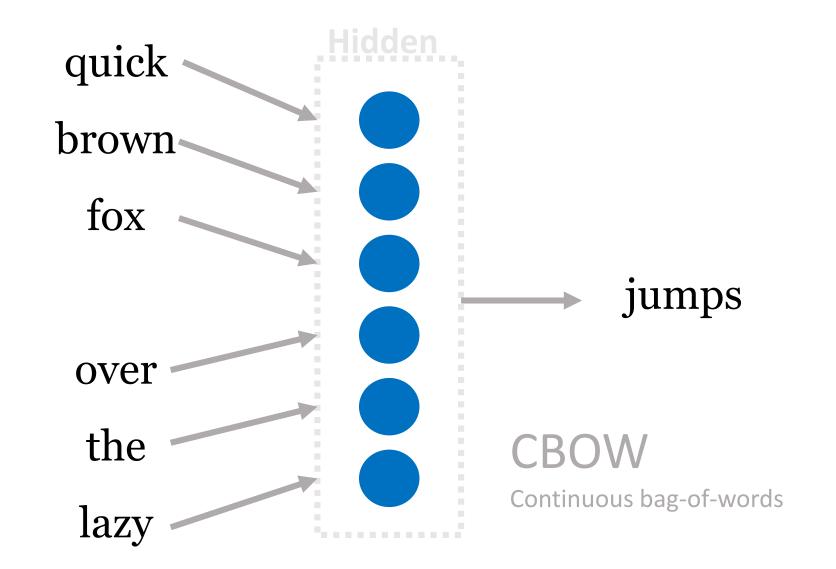
Word2Vec

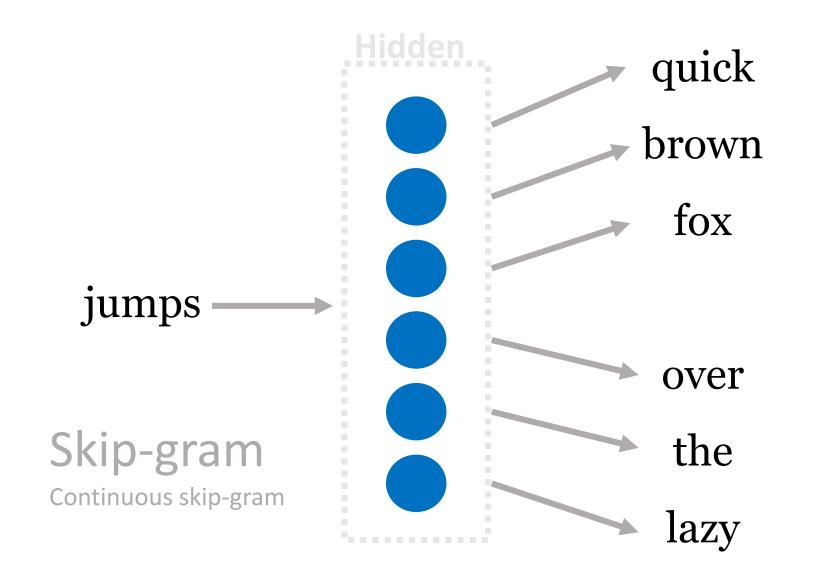


"The quick brown fox ____? over the lazy dog"



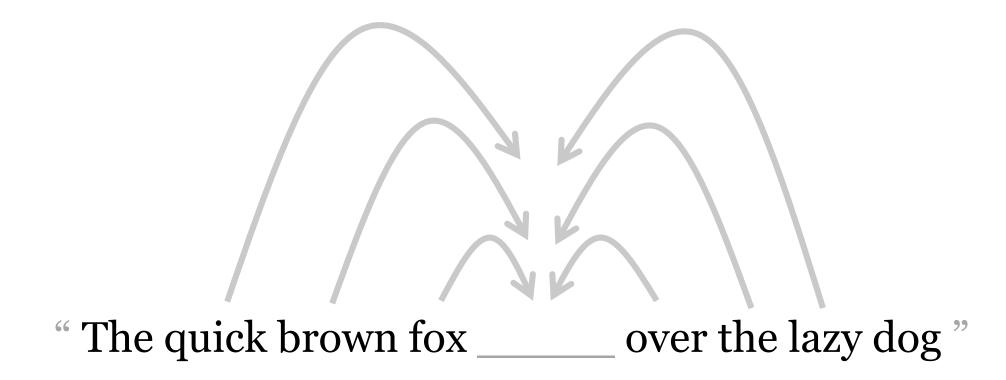






Dimension wood or distation

```
0.0010
0.0000
0.0034
0.2421
0.0000
0.0000
0.0000
```



grammatical, semantical similarity

Word Embedding Choices

- 1. Learnable embedding
- 2. Word2Vec
- 3. GloVe
- 4. FastText

Word2Vec

- Check genism
- •Get the training data on the same page (text8)
- https://radimrehurek.com/gensim/models/word2vec.
 html

Text8

Large text compression benchmark

- First 10⁹ bytes of the English Wikipedia dump on Mar. 3, 2006.
- Remove tags, digits, punctations
- Lower cases
- Leaving a-z, unrepeated spaces
- Truncate first 10⁸ bytes
- About 1700 articles
- http://mattmahoney.net/dc/textdata.html

```
enwik8
                Step
             Original size
100,000,000
96,829,911
             Discard all outside <text...> ... </text>
96,604,864
             Discard #REDIRECT text
96,210,439
             Discard XML tags (<text...> and </text>)
95,287,203
             URL-decode < &gt; and &amp; to < > and &
95,087,290
             Remove <ref> ... </ref> (citations)
93,645,338
             Remove other XHTML tags
             Replace [http:... anchor text] with [anchor text]
91,399,021
90,868,662
             Replace [[Image:...|thumb|left/right|NNNpx|caption]] with caption
             Replace [[category:text|title]] with [[text]]
90,770,617
88,385,654
             Remove [[language:link]] (links to same page in other languages)
             Replace [[Wiki link|anchor text]] with [[anchor text]]
85,443,983
83,420,173
             Remove {{...}} (icons and special symbols)
80,943,967
             Remove { ... } (tables)
77,732,609
                    [ and ]
             Remove
75,053,443
             Replace &...; with space (URL-encoded chars)
             Convert to lower case, replace all sequences not in a-z,0-9 with a single space
70,007,945
74,090,640
             Spell digits, leaving a-z and unrepeated spaces
```

 $O = E \times T \times Q$

- O Training complexity
- E Training epochs
 - T Word count in corpus
 - *Q* <*Model related factor>*

$O = E \times T \times Q$

- O Training complexity
- E Training epochs
- T Word count in corpus
- *Q* <*Model related factor>*

CBOW

Continuous bag-of-words

N Window size for input

D Dimension of P layer

V Size of vocabulary

$$Q = N \times D + D \times log_2(V)$$

$$O = E \times T \times Q$$

- O Training complexity
- E Training epochs
- T Word count in corpus
- *Q* <*Model related factor>*

Skip-gram

Continuous skip-gram

$$Q = C \times (D + D \times log_2(V))$$

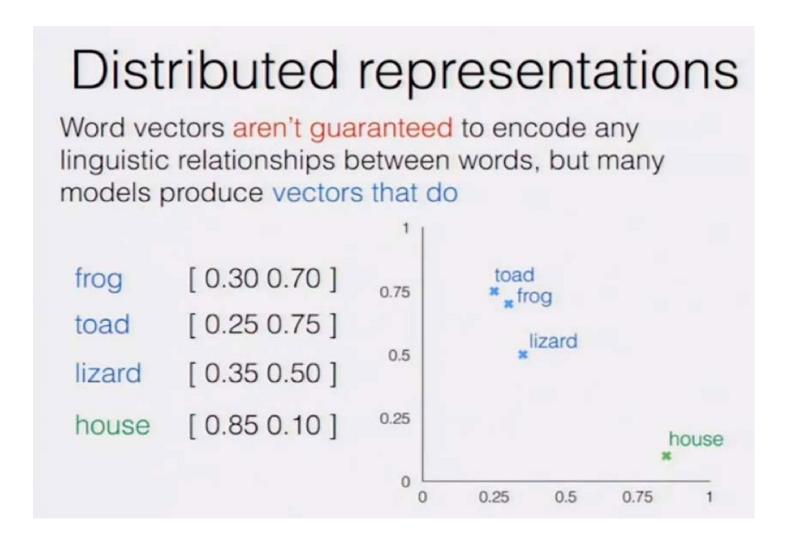
Word count

• Try Wikipedia Statistics page https://en.wikipedia.org/wiki/Wikipedia:Size_in_volumes

Vocabulary size

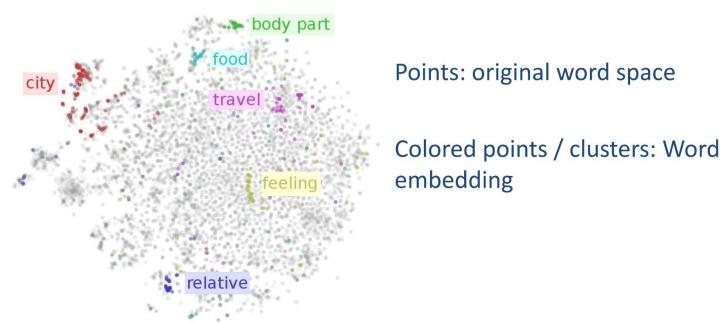
- Try Heaps' law
- https://en.wikipedia.org/wiki/Heaps%27_law

Problem?



Example

Any technique mapping a word (or phrase) from it's original high-dimensional input space (the body of all words) to a lower-dimensional numerical vector space - so one *embeds* the word in a different space

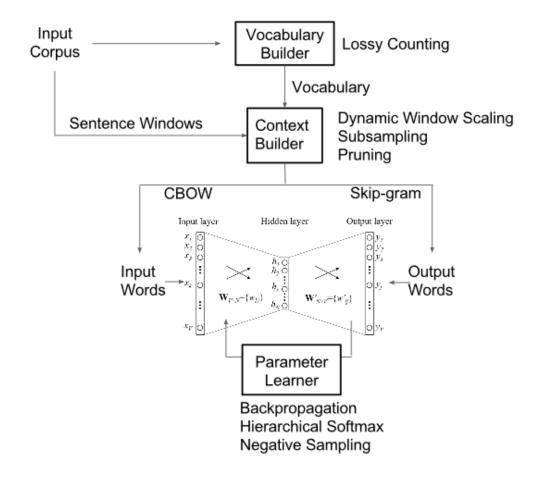


Source: http://sebastianruder.com/content/images/2016/04/word_embeddings_colah.png

Word Representations

Traditional Method - Bag of Words Model	Word Embeddings
 Uses one hot encoding Each word in the vocabulary is represented by one bit position in a 	 Stores each word in as a point in space, where it is represented by a vector of fixed number of dimensions (generally 300)
 For example, if we have a vocabulary of 10000 words, and "Hello" is the 4th word in the dictionary, it would be represented by: 0 0 0 1 0 0 0 0 0 0 	 Unsupervised, built just by reading huge corpus For example, "Hello" might be represented as: [0.4, -0.11, 0.55, 0.3 0.1, 0.02]
Context information is not utilized	Dimensions are basically projections along different axes, more of a mathematical concept.

Architecture





To compare pieces of text

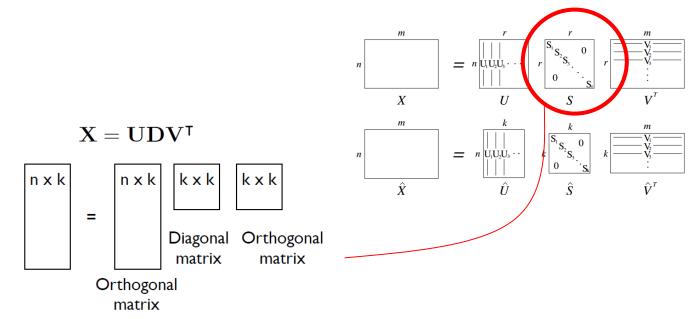
- We need effective representation of
 - Words
 - Sentences
 - Text
- Approach 1: Use existing thesauri or ontologies like WordNet and Snomed CT (for medical).

Drawbacks:

- Manual
- Not context specific
- Approach 2: Use co-occurrences for word similarity. Drawbacks:
 - Quadratic space needed
 - Relative position and order of words not considered

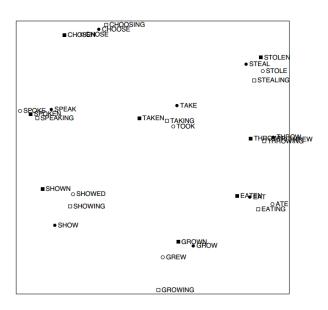
Approach 3: low dimensional vectors

- Store only "important" information in fixed, low dimensional vector.
- Singular Value Decomposition (SVD) on co-occurrence matrix
 - \hat{X} is the best rank *k* approximation to *X*, in terms of least squares
 - Motel = [0.286, 0.792, -0.177, -0.107, 0.109, -0.542, 0.349, 0.271]
- m = n = size of vocabulary
- \hat{S} is the same matrix as S except that it contains only the top largest singular values



Example of Approach 3: low dimensional vectors

• An Improved Model of Semantic Similarity Based on Lexical Co-Occurrence [Rohde et al. 2005]



Problems with SVD

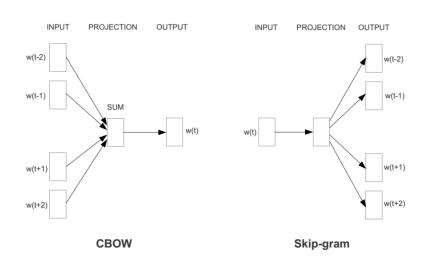
- Computational cost scales quadratically for n x m matrix: O(mn²) flops (when n<m)
- Hard to incorporate new words or documents
- Does not consider order of words

word2vec approach to represent the meaning of word

- Represent each word with a low-dimensional vector
- Word similarity = vector similarity
- Key idea: Predict surrounding words of every word
- Faster and can easily incorporate a new sentence/document or add a word to the vocabulary

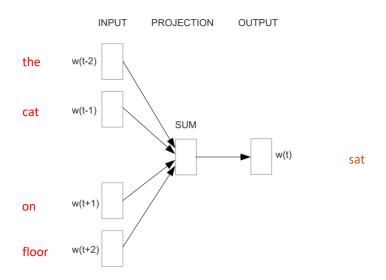
Represent the meaning of word – word2vec

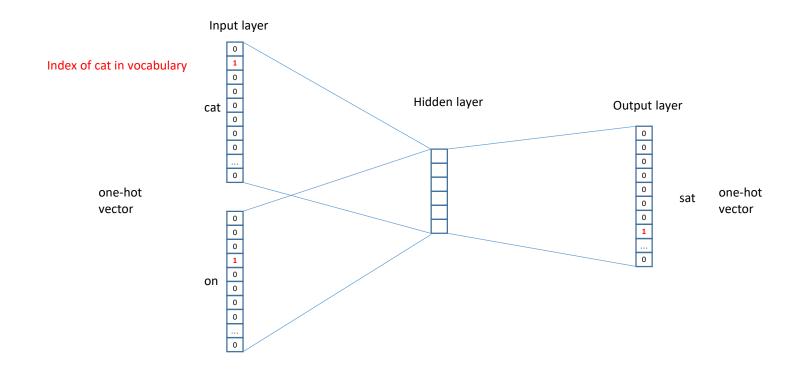
- 2 basic neural network models:
 - Continuous Bag of Word (CBOW): use a window of word to predict the middle word
 - Skip-gram (SG): use a word to predict the surrounding ones in window.

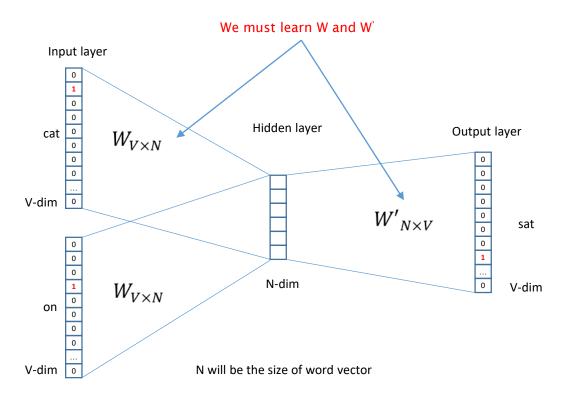


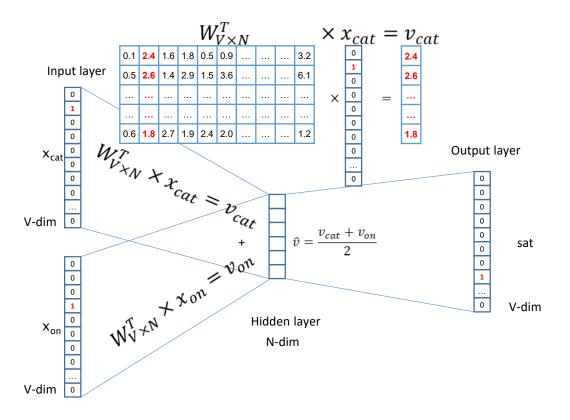
Word2vec – Continuous Bag of Word

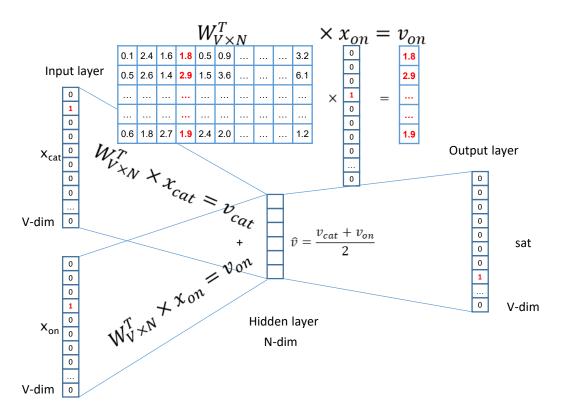
- E.g. "The cat sat on floor"
 - Window size = 2

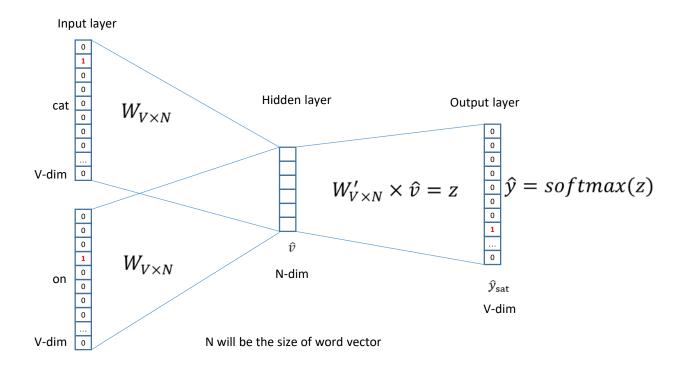


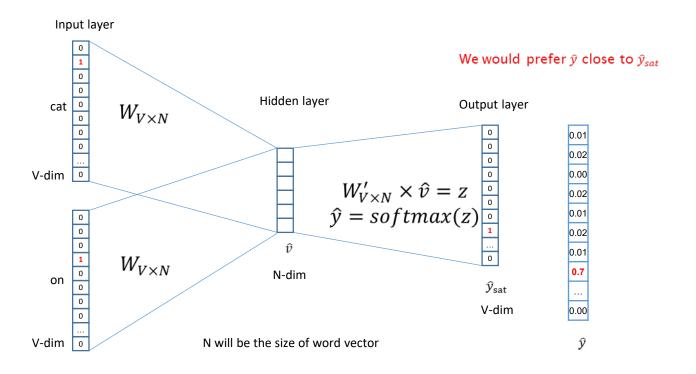


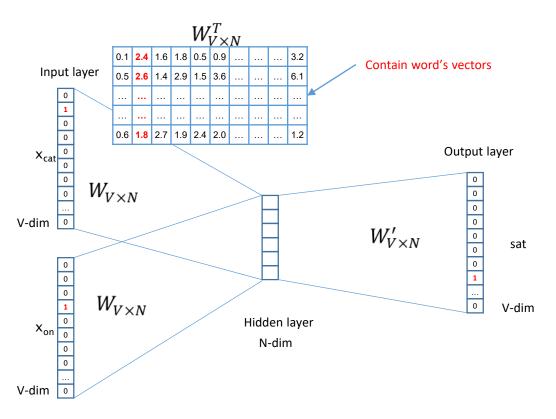










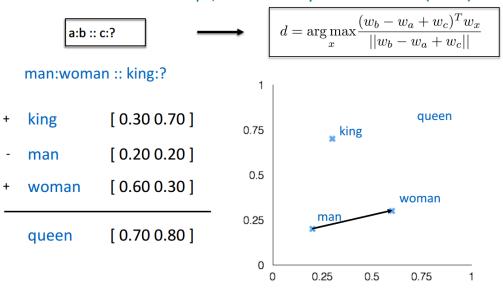


We can consider either W or W' as the word's representation. Or even take the average.

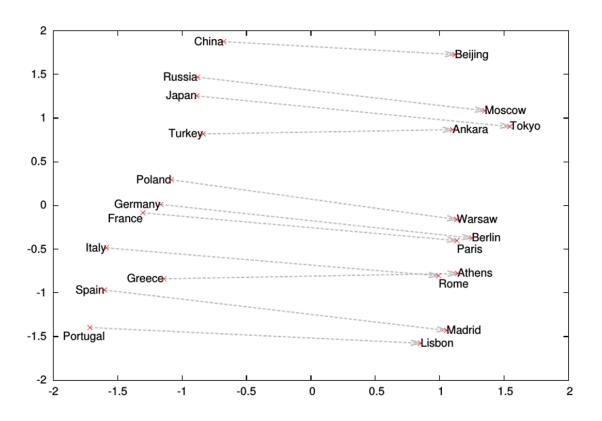
Some interesting results

Word Analogies

Test for linear relationships, examined by Mikolov et al. (2014)

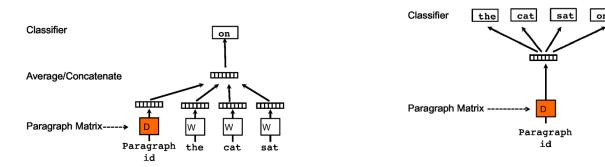


Word analogies



Represent the meaning of sentence/text

- Simple approach: take avg of the word2vecs of its words
- Another approach: Paragraph vector (2014, Quoc Le, Mikolov)
 - Extend word2vec to text level
 - Also two models: add paragraph vector as the input



Applications

- Word Similarity: Edit Distance, WordNet, Porter's Stemmer, Lemmatization using dictionaries
- Search, e.g., query expansion
- Machine Translation
- Part-of-Speech and Named Entity Recognition
- Relation extraction
- Sentiment analysis
- Semantic Analysis of Documents
- Clustering