# Text Analytics

## Introduction to word2vec

By

Kathirmani Sukumar

# Classic Vector Representations

**Document Term Matrix**

|          | $T_1$ | $T_2$ | $T_3$ | $T_4$ | ... | $T_n$ |
|----------|-------|-------|-------|-------|-----|-------|
| $D_1$    | 5     | 0     | ...   | ...   | ... | 0     |
| $D_2$    | 0     | 1     | ...   | ...   | ... | 0     |
| $D_3$    | 0     | 1     | ...   | ...   | ... | 0     |
| $D_4$    | 0     | 1     | ...   | ...   | ... | 1     |
| ...      |       |       | ...   | ...   | ... | 1     |
| $D_m$    | 1     | 2     | ...   | ...   | ... | 10    |

Vector representation of $n^{th}$ term

Vector representation of $m^{th}$ document

# Limitations

▸ Sparse matrix
  ▸ Very less non-zero values. Mostly 90%-95% values are zero
▸ High dimension
  ▸ For each word vector dimension is equal the number of documents
▸ Weak relationship between terms
  ▸ Context between words are loosely represented

| | $T_1$ | $T_2$ | $T_3$ | $T_4$ | ... | $T_n$ |
|---|---|---|---|---|---|---|
| $D_1$ | 5 | 2 | ... | ... | ... | 0 |
| $D_2$ | 2 | 1 | ... | ... | ... | 0 |
| $D_3$ | 0 | 1 | ... | ... | ... | 0 |
| $D_4$ | 1 | 1 | ... | ... | ... | 1 |
| ... | | | ... | ... | ... | 1 |
| $D_m$ | 1 | 2 | ... | ... | ... | 10 |

T1 & T2 has appeared together across many documents. But did they appear next to each other??

# Word Embeddings

▶ Word embeddings is a collective name used for those techniques where words are translated in to dense low dimensional vectors instead of sparse high dimensional vectors

▶ These techniques are usually driven by neural network based model compared to traditional frequency based models (like LSA)

▶ Word embeddings techniques
  ▶ Word2vec (developed by Google)
  ▶ Global Vectors for Word Representation - GloVe (developed by Standford)
  ▶ FastText

# word2vec

▶ Word2vec is a computationally-efficient predictive model for learning word embeddings from raw text.



Input: Raw Text

Source: orielly.com

Output: D – Dimensional Dense Vector for each word

# Google Pre-Trained Model

▸ Difficult to train our own model for large text. Need high end machines

▸ Input: Google has used Google News data set (about 100 billion words) to train a word2vec model.

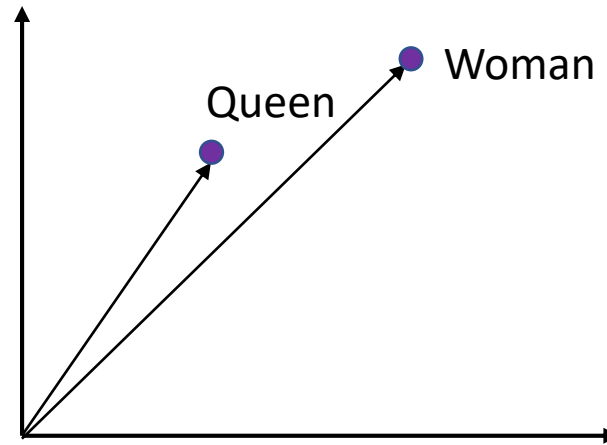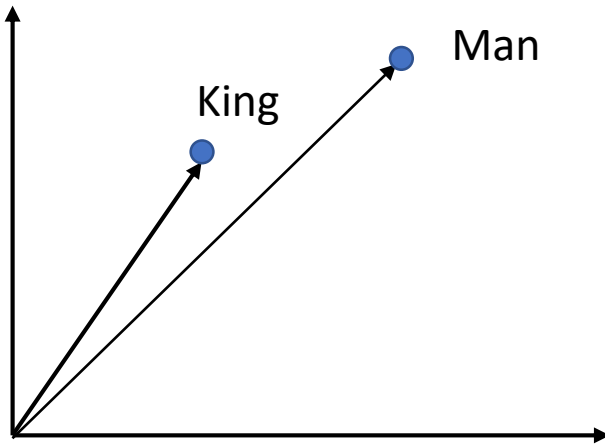▸ Output: Each word is represented using 300 dimension

```python
print(model.get_vector('computer'))
```

```
[ 1.07421875e-01 -2.01171875e-01  1.23046875e-01  2.11914062e-01
 -9.13085938e-02  2.16796875e-01 -1.31835938e-01  8.30078125e-02
  2.02148438e-01  4.78515625e-02  3.66210938e-02 -2.45361328e-02
  2.39257812e-02 -1.60156250e-01 -2.61230469e-02  9.71679688e-02
 -6.34765625e-02  1.84570312e-01  1.70898438e-01 -1.63085938e-01
 -1.09375000e-01  1.49414062e-01 -4.65393066e-04  9.61914062e-02
  1.68945312e-01  2.60925293e-03  8.93554688e-02  6.49414062e-02
  3.56445312e-02 -6.93359375e-02 -1.46484375e-01 -1.21093750e-01
 -2.27539062e-01  2.45361328e-02 -1.24511719e-01 -3.18359375e-01
 -2.20703125e-01  1.30859375e-01  3.66210938e-02 -3.63769531e-02
 -1.13281250e-01  1.95312500e-01  9.76562500e-02  1.26953125e-01
  6.59179688e-02  6.93359375e-02  1.02539062e-02  1.75781250e-01
 -1.68945312e-01  1.21307373e-03 -2.98828125e-01 -1.15234375e-01
  5.66406250e-02 -1.77734375e-01 -2.08984375e-01  1.76757812e-01
  2.38037109e-02 -2.57812500e-01 -4.46777344e-02  1.88476562e-01
  5.51757812e-02  5.02929688e-02 -1.06933594e-01  1.89453125e-01
 -1.16210938e-01  8.49609375e-02 -1.71875000e-01  2.45117188e-01
 -1.73828125e-01 -8.30078125e-03  4.56542969e-02 -1.61132812e-02
  1.86523438e-01 -6.05468750e-02 -4.17480469e-02  1.82617188e-01
```
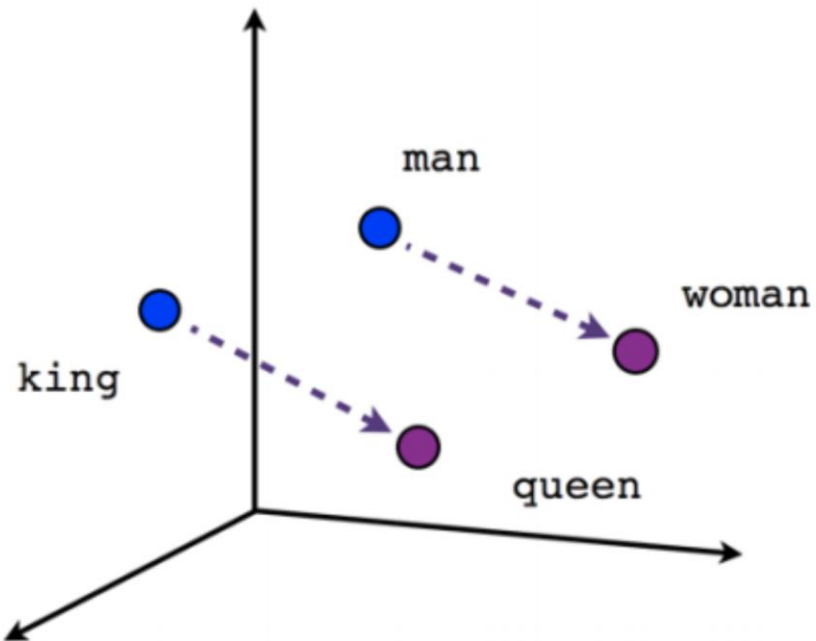
Link to download the model: https://drive.google.com/file/d/0B7XkCwpI5KDYNlNUTTlSS21pQmM/edit

# Characteristic

▶ Similar words have similar vector representation

# Characteristic

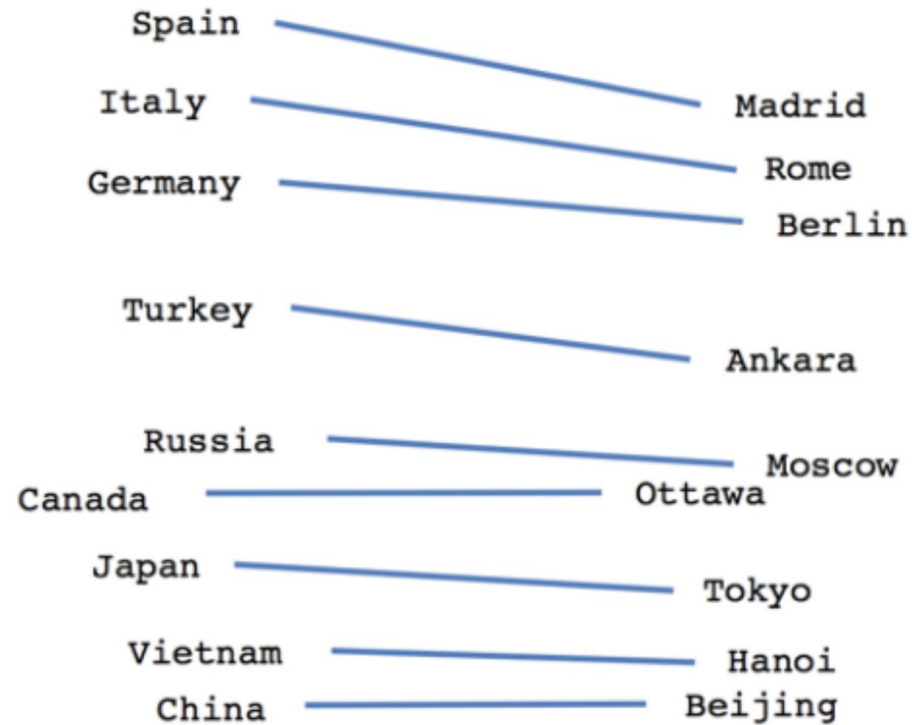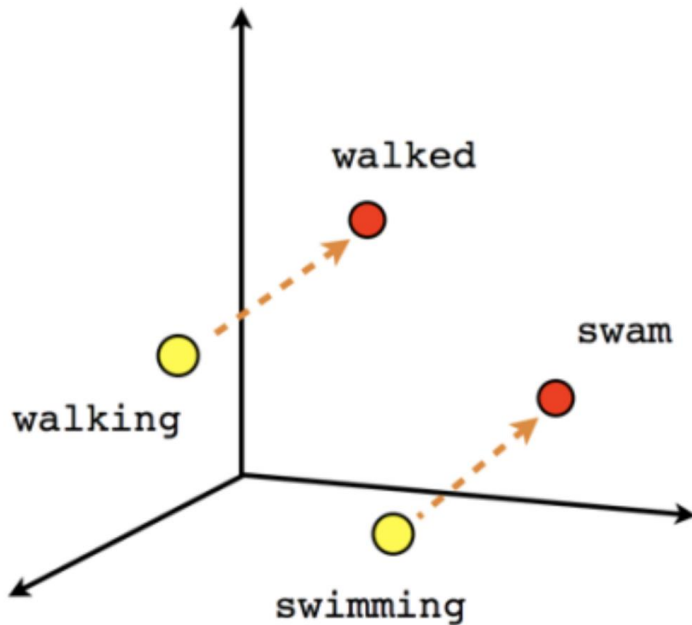▶ Similar words have similar vector representation



**King – Man + Woman = Queen**

**King – Man = Queen - Woman**

Source: https://www.tensorflow.org/tutorials/representation/word2vec

# Characteristic

▸ Similar words have similar vector representation



Source: https://www.tensorflow.org/tutorials/representation/word2vec

# Applications

▸ Word similarity
  ▸ Thesaurus: Given a word, we can identify its related words
▸ Stemming
▸ Parts of Speech
▸ Named Entity Recognition
▸ Sentiment Analysis
▸ Recommendation Engines

# THANK YOU

**manipal PROlearn**