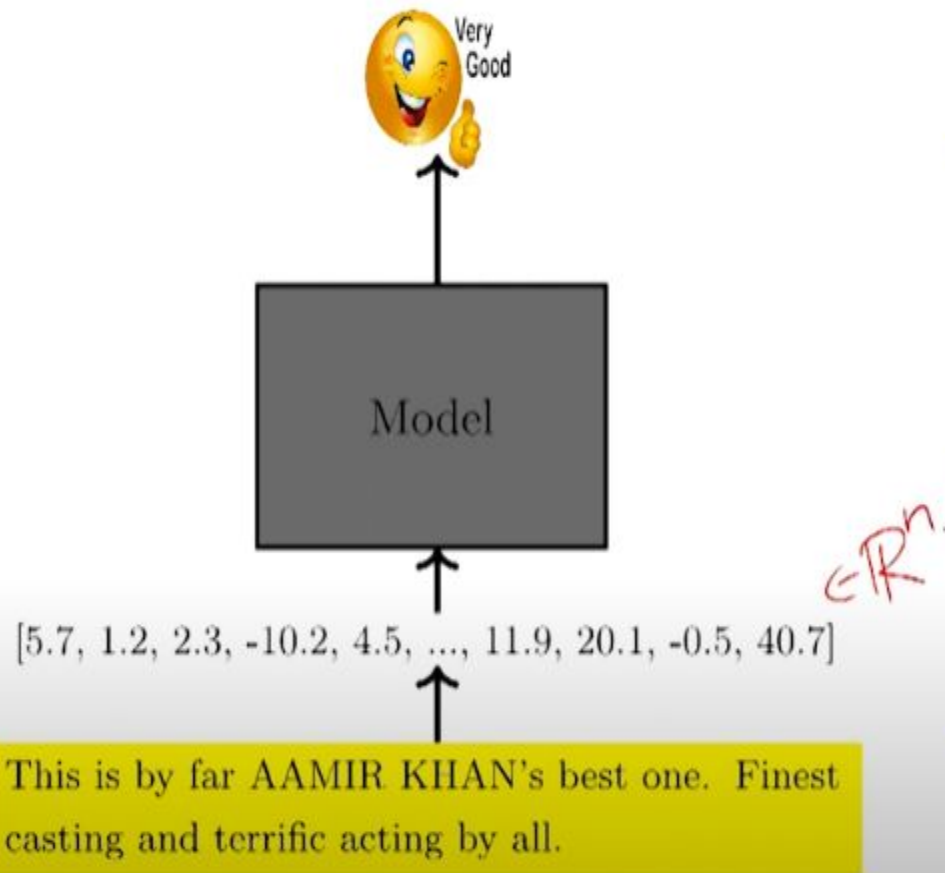# Word to Vector Representation

- Let us start with a very simple motivation for why we are interested in vectorial representations of words

- Suppose we are given an input stream of words (sentence, document, etc.) and we are interested in learning some function of it (say, $\hat{y} = sentiments(words)$)

- Say, we employ a machine learning algorithm (some mathematical model) for learning such a function ($\hat{y} = f(\mathbf{x})$)

- We first need a way of converting the input stream (or each word in the stream) to a vector $\mathbf{x}$ (a mathematical quantity)

Model

$\in \mathbb{R}^n$

$[5.7, 1.2, 2.3, -10.2, 4.5, ..., 11.9, 20.1, -0.5, 40.7]$

This is by far AAMIR KHAN's best one. Finest casting and terrific acting by all.

Very Good

## Corpus:

- Human machine interface for computer applications

- User opinion of computer system response time

- User interface management system

- System engineering for improved response time

**V** = [human,machine, interface, for, computer, applications, user, opinion, of, system, response, time, interface, management, engineering, improved]

- Given a corpus, consider the set $V$ of all unique words across all input streams (*i.e.*, all sentences or documents)

- $V$ is called the **vocabulary** of the corpus (*i.e.*, all sentences or documents)

- We need a representation for every word in $V$

- One very simple way of doing this is to use one-hot vectors of size $|V|$

cat:

| 0 | 0 | 0 | 0 | 0 | 1 | 0 |
|---|---|---|---|---|---|---|

dog:

| 0 | 1 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|

truck:

| 0 | 0 | 0 | 1 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|

$$euclid\_dist(\mathbf{cat}, \mathbf{dog}) = \sqrt{2}$$
$$euclid\_dist(\mathbf{dog}, \mathbf{truck}) = \sqrt{2}$$
$$cosine\_sim(\mathbf{cat}, \mathbf{dog}) = 0$$
$$cosine\_sim(\mathbf{dog}, \mathbf{truck}) = 0$$

**Problems:**

- $V$ tends to be very large (for example, 50K for PTB, 13M for Google 1T corpus)
- These representations do not capture any notion of similarity
- Ideally, we would want the representations of cat and dog (both domestic animals) to be closer to each other than the representations of cat and truck
- However, with 1-hot representations, the Euclidean distance between **any two words** in the vocabulary in $\sqrt{2}$

And the cosine similarity between **any two words** in the vocabulary is 0

- *You shall know a word by the company it keeps - Firth, J. R. 1957:11*
- Distributional similarity based representations
- This leads us to the idea of co-occurrence matrix

A bank is a **financial** institution that accepts **deposits** from the public and creates **credit**.

The idea is to use the accompanying words (financial, deposits, credit) to represent bank

## Corpus:

- Human machine interface for computer applications

- User opinion of computer system response time

- User interface management system

- System engineering for improved response time

- A co-occurrence matrix is a **terms** × **terms** matrix which captures the number of times a term appears in the context of another term

## Corpus:

- Human machine interface for computer applications
- User opinion of computer system response time
- User interface management system
- System engineering for improved response time

| | human | machine | system | for | ... | user |
|---|---|---|---|---|---|---|
| human | 0 | 1 | 0 | 1 | ... | 0 |
| machine | 1 | 0 | 0 | 1 | ... | 0 |
| system | 0 | 0 | 0 | 1 | ... | 2 |
| for | 1 | 1 | 1 | 0 | ... | 0 |
| . | . | . | . | . | . | . |
| . | . | . | . | . | . | . |
| . | . | . | . | . | . | . |
| user | 0 | 0 | 2 | 0 | ... | 0 |

Co-occurence Matrix

- A co-occurrence matrix is a **terms × terms** matrix which captures the number of times a term appears in the context of another term
- The context is defined as a window of $k$ words around the terms
- Let us build a co-occurrence matrix for this toy corpus with $k = 2$
- This is also known as a **word × context** matrix
- You could choose the set of **words** and **contexts** to be same or different
- Each row (column) of the co-occurrence matrix gives a vectorial representation of the corresponding word (context)

|        | human | machine | system | for | ... | user |
|--------|-------|---------|--------|-----|-----|------|
| human  | 0     | 1       | 0      | x   | ... | 0    |
| machine| 1     | 0       | 0      | x   | ... | 0    |
| system | 0     | 0       | 0      | x   | ... | 2    |
| for    | x     | x       | x      | x   | ... | x    |
| .      | .     | .       | .      | .   | .   | .    |
| user   | 0     | 0       | 2      | x   | ... | 0    |

## Some (fixable) problems

- Stop words (a, the, for, etc.) are very frequent → these counts will be very high

- Solution 1: Ignore very frequent words

- Solution 2: Use a threshold t (say, t = 100)

$$X_{ij} = min(count(w_i, c_j), t),$$

where $w$ is word and $c$ is context.

## Some (severe) problems

- Very high dimensional ($|V|$)

- Very sparse

- Grows with the size of the vocabulary

- **Solution:** Use dimensionality reduction (SVD)

# Continuous Bag of Words (CBOW)

- The methods that we have seen so far are called **count based models** because they use the co-occurrence counts of words
- We will now see methods which directly **learn** word representations (these are called **(direct) prediction based models**)

## What is Word2Vec ?

- A two layer neural network to generate word embeddings given a text corpus.

- Word Embeddings – Mapping of words in a vector space.

| Man | Women |
|---|---|
| 0.52 | 0.73 |
| 0.76 | 0.89 |
| 1.21 | -1.67 |
| 0.22 | 1.32 |
| -1.36 | 0.36 |
| 0.49 | -1.49 |
| -3.69 | 2.71 |
| -0.07 | 0.05 |

# Why Word2vec?

- Preserves relationship between words.
- Deals with addition of new words in the vocabulary.
- Better results in lots of deep learning applications.

# Working of word2Vec

- The word2vec objective function causes the words that occur in similar contexts to have similar embeddings.

Example: The **kid** said he would grow up to be superman.

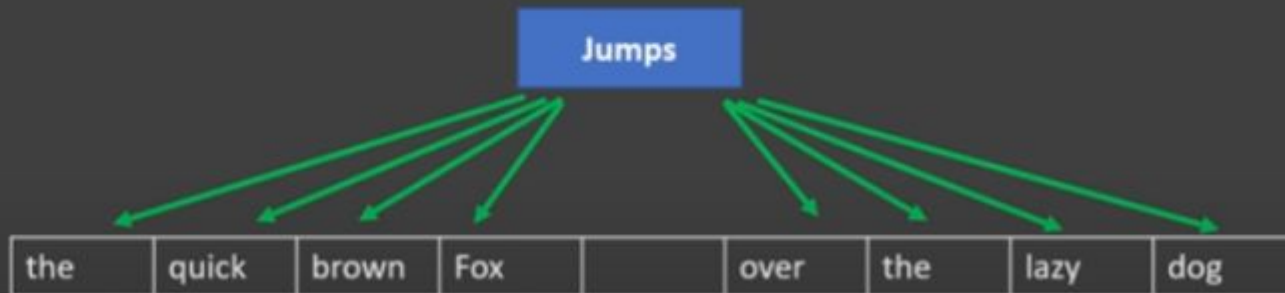The **child** said he would grow up to be superman.

# CBOW

- Predict the target word from the context.

| the | quick | brown | fox | | over | the | lazy | dog |

**Jumps**

# Skip Gram

- Predict the context words from target.

**Jumps**

| the | quick | brown | Fox | | over | the | lazy | dog |

1. Take a fake problem
2. Solve it using neural network
3. You get word embeddings as a side effect

fake problem: fill in a missing word in a sentence

There lived a king called Ashoka in India. After Kalinga battle, he converted to Buddhism. This **mighty king ordered his ministers** to put together a peaceful treaty with their neighboring kingdoms. The **emperor ordered his ministers to** also build stupa, a monument with Buddha's teachings.

Fake problem

_____ ordered his ministers

_____ ordered his ministers

There lived a king called Ashoka in India. After Kalinga battle, he converted to Buddhism. This mighty king ordered his ministers to put together a peaceful treaty with their neighboring kingdoms. The emperor ordered his ministers to also build stupa, a monument with Buddha's teachings.

**Fake problem**

king ordered his ministers

emperor ordered his ministers

**Side effect**

$$king \begin{bmatrix} 0.7 \\ 0.4 \\ 1.2 \\ 3.8 \end{bmatrix} \quad emperor \begin{bmatrix} 0.7 \\ 0.5 \\ 1.2 \\ 3.8 \end{bmatrix} \quad king \sim emperor$$

CODE BASICS

eating _____ is very healthy          table, angry, truck, apple, pizza, walnut

NASA launched _____ last month          table, angry, truck, rocket, apple, pizza

There lived a king called Ashoka in India. After Kalinga battle, he converted to Buddhism. This mighty king ordered his ministers to put together a peaceful treaty with their neighboring kingdoms. The emperor ordered his ministers to also build stupa, a monument with Buddha's teachings.

Training samples

lived, a → There

a, king → lived

ordered, his → king

ordered, his → emperor

emperor **ordered**
**his**

ordered

| | |
|---|---|
| Ashoka | 0 |
| emperor | 0 |
| his | 0 |
| king | 0 |
| ordered | 1 |
| ... | ... |
| zone | 0 |

his

| | |
|---|---|
| Ashoka | 0 |
| emperor | 0 |
| his | 1 |
| king | 0 |
| ordered | 0 |
| ... | ... |
| zone | 0 |

$\hat{y}$

| 0.07 | | Ashoka |
| 0.1 | | emperor |
| 0.4 | | his |
| 0.23 | | king |
| 0.00 | | ordered |
| ... | | ... |
| 0.09 | | zone |

$y$

| 0 |
| 1 |
| 0 |
| 0 |
| 0 |
| ... |
| 0 |

$\Sigma \sigma$
$\Sigma \sigma$
$\Sigma \sigma$
$\Sigma \sigma$

# CBOW - Working

Hope can set you free.



$V_{5 \times 1}$, one hot vector of "Hope"

$W_{3 \times 5}$

3 nodes in hidden layer

$W'_{5 \times 3}$

Compare and Update weights

Actual Target

$V_{5 \times 1}$, one hot vector of "Set"

$W_{3 \times 5}$

$V_{5 \times 1}$, predicted one hot vector of "Can"

| w00 | w01 | w02 | w03 | w04 |
|-----|-----|-----|-----|-----|
| w10 | w11 | w12 | w13 | w14 |
| w20 | w21 | w22 | w23 | w24 |

$W_{3 \times 5}$

# CBOW - Working

Hope can **set you free**

$V_{5 \times 1}$, one hot vector of "Hope"

$W_{3 \times 5}$

3 nodes in hidden layer

$W'_{5 \times 3}$

Compare and Update weights

Actual Target

$V_{5 \times 1}$, one hot vector of "Set"

$W_{3 \times 5}$

$V_{5 \times 1}$, predicted one hot vector of "Can"

| w00 | w01 | w02 | w03 | w04 |
|-----|-----|-----|-----|-----|
| w10 | w11 | w12 | w13 | w14 |
| w20 | w21 | w22 | w23 | w24 |

$W_{3 \times 5}$

# Getting word embeddings

One Hot vector of words $V_{5 \times 1}$

Weights after training
$W_{3 \times 5}$

| w00 | w01 | w02 | w03 | w04 |
|-----|-----|-----|-----|-----|
| w10 | w11 | w12 | w13 | w14 |
| w20 | w21 | w22 | w23 | w24 |

| 1 | | 0 | | 0 | | 0 | | 0 |
|---|---|---|---|---|---|---|---|---|
| 0 | | 1 | | 0 | | 0 | | 0 |
| 0 | | 0 | | 1 | | 0 | | 0 |
| 0 | | 0 | | 0 | | 1 | | 0 |
| 0 | | 0 | | 0 | | 0 | | 1 |

Hope    can    set    you    free

Word Vector for hope $= W_{3 \times 5} \ X \ V_{5 \times 1}$

$V_{3 \times 1}$

| w00 | w01 | w02 | w03 | w04 |
|-----|-----|-----|-----|-----|
| w10 | w11 | w12 | w13 | w14 |
| w20 | w21 | w22 | w23 | w24 |

X

| 1 |
|---|
| 0 |
| 0 |
| 0 |
| 0 |

=

| w00 |
|-----|
| w10 |
| w20 |

# Improving the accuracy

- Choice of Model architecture (CBOW / Skipgram
  - Large Corpus, higher dimensions, slower– Skipgram
  - Small Corpus, Faster - CBOW
- Increasing the training dataset.
- Increasing the vector dimensions
- Increasing the windows size.

# Skip Gram Model

Skip Gram - Working

Hope can set you free.

# Getting word embeddings

One Hot vector of words     $V_{5 \times 1}$

Weights after training

$W_{3 \times 5}$

| w00 | w01 | w02 | w03 | w04 |
|-----|-----|-----|-----|-----|
| w10 | w11 | w12 | w13 | w14 |
| w20 | w21 | w22 | w23 | w24 |

| | | | | |
|---|---|---|---|---|
| 1 | 0 | 0 | 0 | 0 |
| 0 | 1 | 0 | 0 | 0 |
| 0 | 0 | 1 | 0 | 0 |
| 0 | 0 | 0 | 1 | 0 |
| 0 | 0 | 0 | 0 | 1 |
| Hope | can | set | you | free |

Word Vector for hope $= W_{3 \times 5} \ X \ V_{5 \times 1}$

$V_{3 \times 1}$

| w00 | w01 | w02 | w03 | w04 |
|-----|-----|-----|-----|-----|
| w10 | w11 | w12 | w13 | w14 |
| w20 | w21 | w22 | w23 | w24 |

$X$

| |
|---|
| 1 |
| 0 |
| 0 |
| 0 |
| 0 |

$=$

| |
|---|
| w00 |
| w10 |
| w20 |