# Jay Alammar

# The Illustrated Word2Vec

→ Word2vec
  ↳ efficiently create
    embeddings
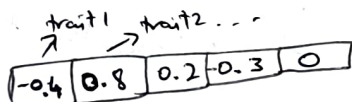
Example :

→ personality trait
  any person's personality can be
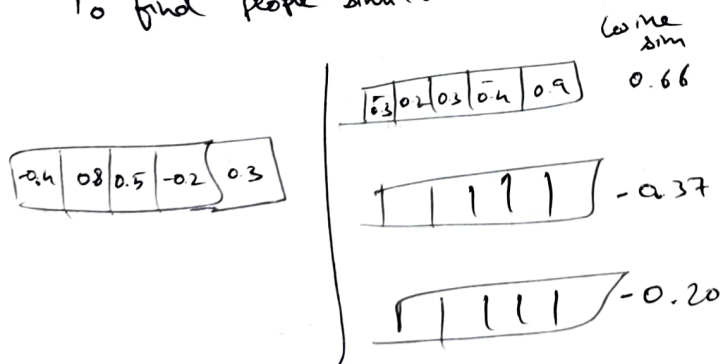  expressed in 5 numbers
  features are   i) extraversion
                 ii) negative emotionally
                 :

  You can express a person in
  a  5 dimensional  vector space
        trait1  trait2 ...
        ↗      ↗
  | -0.4 | 0.8 | 0.2 | -0.3 | 0 |

to detect 2 persons with similar
personality   we   plot them in a
5 dimensional vector space
_ and use algorithms to calculate
how close each person is
for example using (cosine similarity)

To find people similar to   you,

| 0.4 | 08 | 0.5 | -0.2 | 0.3 |

| -0.3 | 0.2 | 0.3 | -0.4 | 0.9 |   Cosine sim  0.66

| 1 | 1 | 1 | 1 | 1 |   -0.37

| 1 | 1 | 1 | 1 | 1 |   -0.20

0.66 % The guy most similar to
         you

---

## Word embeddings

GloVe vector trained on
wikipedia has 50 numbers to
define a word

we don't need what those
50 numbers represents

vector embedding of

king + woman  - man

      gives   queen
              throne
              prince
              :

now looking at training the
word embedding model

Language model training

word1 →
word2 →   embed  →  Project → Logits for
  :   →   ↳ predict↗      each
                           word in
                           vocab

data :
  we  get a paragraph
  run a  window over the text
  and  that text will be
  input  to the model when

     X = text - window
     y = window

# Skipgram dataset

→ in a text blob, we get a sliding window say with 5 words as window size, we take the middle text as ~~text~~ ✗ and the remaining 4 as y for the model, in this text the skipgram will look something like

| | | | |
|---|---|---|---|
| text | in | text | in |
| text | a | text | a |
| text | blob | text | blob |
| text | we | text | we |
| ~~text~~ | | blob | a |
| blob | a | blob | text |
| blob | text | blob | we |
| | | blob | get |

## Training process

not → | untrained model | → logits

Actual target          model prediction  Error

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

model

| |
|---|
| 0.1 |
| 0.5 |
| : |
| 0.02 |
| : |

| |
|---|
| : |
| : |
| : |
| : |

This model is the embeddings matrix

---

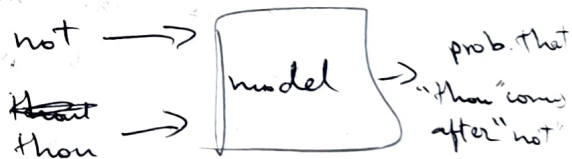not → | untrained model to predict neighbour word | → logits

1) lookup embeddings    2) Calculate precision    3) Project to output

The third step is computationally expensive, to do this we are change our task into 2 steps

1) Generate high quality word embeddings

2) Use this high quality embed. to train our model

→ change the task objective

not → | model | → thou

to

not ⟶ | model | → prob. that
thou ⟶                "thou" comes
                      after "not"

So we change the dataset to

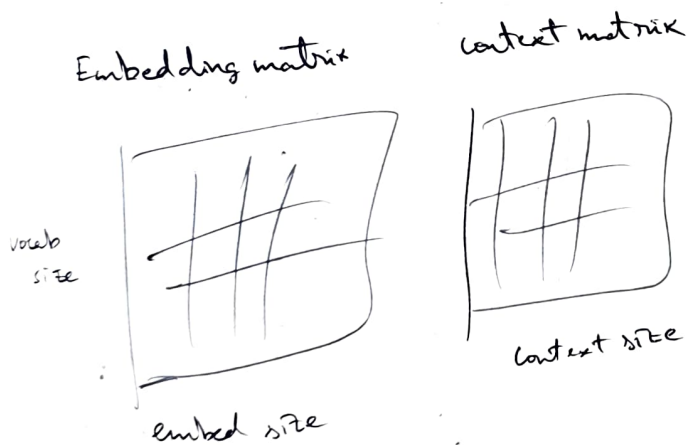| input | | output |
|---|---|---|
| text | in | : |
| text | a | : |
| text | blob | : |
| text | we | : |

but the model will collapse at always predicting "i" so we introduce negative samples by randomly sampling words from the vocabulary

now finally,

## actual <u>word 2 vec</u>

### <u>traing process</u>

define ~~to~~ vocab_size = 10,000
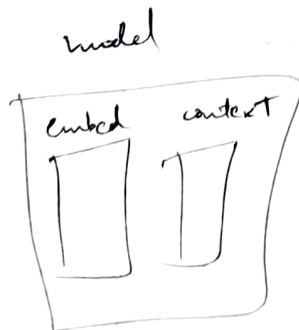
create embedding matrix, context matrix

Embedding matrix          context matrix

vocab size

embed size

context size

Initialize both with random

dataset

| text | in |
| text | a |
| text | blob |
| text | we |
| text | sun |
| text | moon |

| 1 |
| 1 |
| 1 |
| 1 |
| 0 |
| 0 |

model

embed    context

↓        ↓
input    context
word     word

---

~~to~~

## traing step

for input word look in embedding matrix , context word look in context matrix

dot product of

input    opu

Sigmoid (input • output) = model output

calculate error use error to update embed & context table

discard the context table

hyperparameters
→ window size
→ num negative samples

less than 15
⇒ words are interchangeble

>15
⇒ related words

the dataset always contain

1 +ve
n -ve
1 +ve
n -ve
:

and you select n+1 samples per traing step