

- vector semantics is the standard way to represent word meaning in nlp.
- suppose the meaning of the borrowed ongchoi is unknown but we see it in a context:
 - "ongchoi is a delicious sauteed with garlic"
 - "ongchoi leaves with salty sauce"
 - "ongchoi is superb over rice."Because ongchoi appears with words like rice, garlic, delicious, salty, we can infer it is also a leafy green.
computationally, this is done by counting word contexts.

- vector semantics = representing the meaning of words as mathematical vectors each representing the point in a multidimensional space.
- Embeddings = the actual vector representation of a word.

TF:

a term-document matrix has:

- Row_s = word_s
- column_s = document_s
- cell_s = frequency count

The idea is: "Similar documents have similar word distribution."

Since the vocabulary is huge, these vectors are sparse.

■ Term-Term Matrix: ((size |V| × |V|))

Row = ((target word)) column = ((context word))

cells = ((the number of times the target word would happen in some context))

■ "Cosine Similarity" for measuring similarity:

a word or a document is represented by a vector.

$$w = (w_1, w_2, \dots, w_n), v = (v_1, v_2, \dots, v_n) \Rightarrow w \cdot v = \sum_{i=1}^n w_i \cdot v_i$$

What is the problem with Raw Dot Product??

I it depends on vector length

II Frequent words (like 'good', 'the', -) will have unfairly high values.

this means if we consider a $w \cdot v \geq M$, there will be frequent words w_2, v_2 that $w_2 \cdot v_2 > M$. this questions the validity of M as an optimized criteria

→ The idea is to divide $\theta = \text{angle between } w \text{ and } v$ by magnitude $\cos(\theta) = \frac{w \cdot v}{\|w\| \|v\|}$

- What's happening before ??

- we build a cooccurrence matrix ($W \times W$)
- each cell is just raw frequency.

- What is the problem ??

- Raw counts are skewed. Some words appear everywhere
- Common words are not useful for distinguishing
- ~~meaning and context~~
- as a result, similarity matrices get dominated by frequently used uninformative words. which reduces accuracy when comparing meanings.

TF-IDF: make vectors more informative by reducing the weight of common words and increasing weight of rare but discriminative words.

- TF-IDF stands for:

$$\text{Term Frequency} = \text{TF}(w, d) = \frac{\text{count of } w \text{ in } d}{\text{total words in } d}$$

$$\text{Inverse Frequency document} = \text{IDF}(w) = \lg \left(\frac{N}{df(w)} \right)$$

N = total number of documents, $df(w)$ = number of documents that w appeared in

$$\text{TF-IDF}(w, d) = \text{TF}(w, d) \times \text{IDF}(w)$$

- if a word appears a lot in a document, also appears a lot in all documents, it is meaningless stop words.

- We want a way to measure how strongly two words are associated???

- That's what PMI does??

"do these words occur together more than we would expect by chance???"

- $\text{PMI} = \text{pointwise mutual information}$??

$$\text{PMI}(w, c) = \lg \left(\frac{P(w, c)}{P(w)P(c)} \right)$$

$P(w, c)$ = probability that target word w and context-word c occur
 $P(w)$ = word occurring $P(c)$ = context occurring

- if w and c occur more than chance $\Rightarrow \text{PMI}$ is positive
- if w and c occur less than chance $\Rightarrow \text{PMI}$ is negative
- If w and c are independent $\Rightarrow \text{PMI} = 0$

- $\text{PPMI} = \text{positive PMI}$

PMI can be negative, but negative association are often not useful for building semantic vectors.

$$\text{PPMI}(w, c) = \max(\text{PMI}(w, c), 0)$$

$$P(w, c) = \frac{\text{Count}(w, c)}{\text{total cooccurrences}}$$

$$P(w) = \frac{\text{Count}(w)}{\text{total cooccurrences}}$$

- context distribution smoothing

-Instead of using $P(c)$, raise it to the power α (often $\alpha = 0.75$)
 -this makes rare context words less dominant

$$\text{PPMI}_\alpha(w, c) = \max \left(\lg \left(\frac{P(w, c)}{P(w)P(c)^\alpha} \right), 0 \right)$$

• Sparse vectors problems

- very high dimensions - harder to train classifiers
- classifiers will need to many parameters.
- Synonyms are far apart. eg: car have unrelated dimension automobile.

• Dense word embeddings.

- Dense vectors: every number can be non-zero and real (eg: -0.3)
- low-dimensional: only 50-1000 numbers per word.
- no explicit meaning per dimension - but they capture semantic similarity
- words in similar contexts get similar embeddings, even if they never appear together directly.

• Ward-2Vec : two main training algorithms.

- skip-gram with negative sampling ((SGNS))
 - continuous Bag of words ((CBOW))
- pretrained models are available and widely used
produces static embeddings; one fixed vector per word.

• How skip-gram works??

① Suppose we have a corpus "The cat sat on the mat" with a window of size 2, the target-context pairs are:

target	context	
cat	the, sat	⇒ the context is n words before and after target - here n=2
sat	cat, on	
on	the, sat	
the	on, mat	
mat	the	

these pairs are called "positive pairs". Because they happen next to each other.

(cat, the), (cat, sat), (sat, cat), (sat, on),

② Prediction task : given a target word, predict which words are nearby. This is a binary classification task.

- input : (target word vector, also candidate context word embedding)
- predict one for real context words. Of all real words we don't care about this model's prediction, only that while learning to predict, the model adjusts the word vectors.
- but this will be very expensive. the vocabulary might have 100000 words.

③ Negative sampling:

- For real (target, context) pair ((positive ones)), we also make K random (target, random-word-context) pair.
- this model learns to push positive vectors closer, and push negative pairs apart.

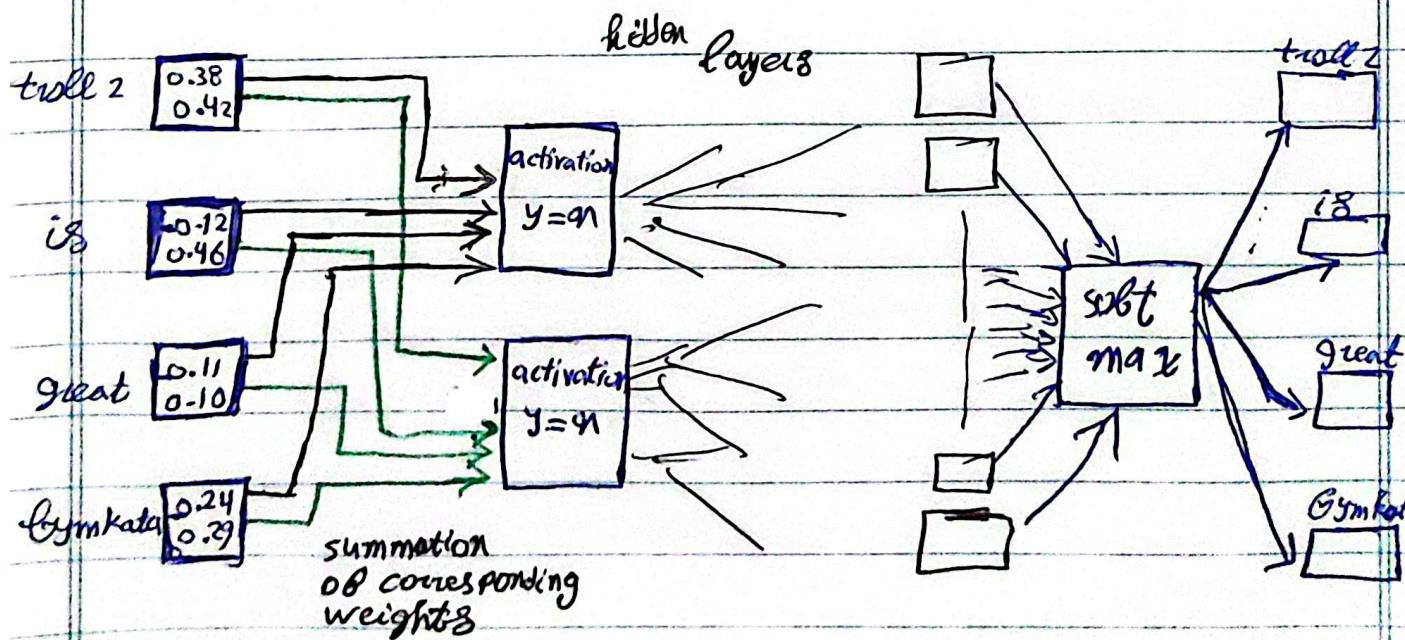
how Word2Vec works (by me)

- how we can get a super simple neural net that figures out what numbers should go with different words. lets imagine we have 2 phrases.

- Troll2 is great
GymRata is great

troll2 and GymRata
are name of movies.

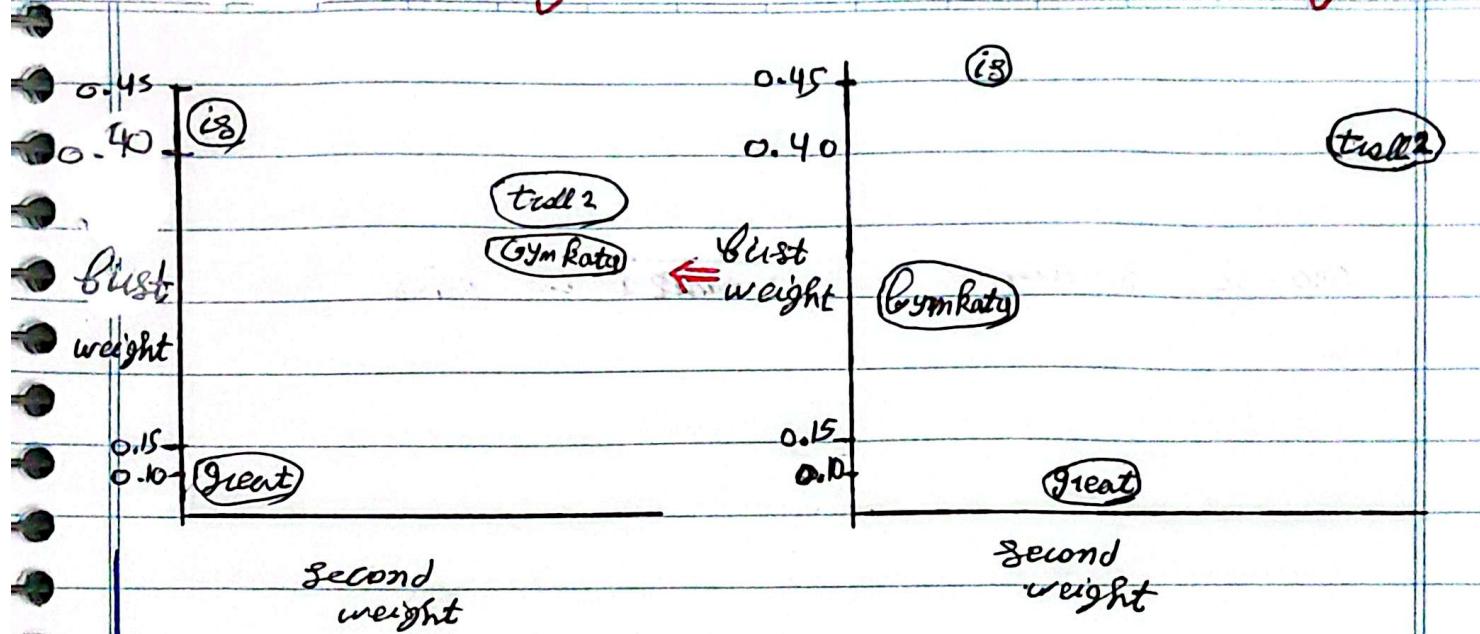
both phrases are about two movies. in order to create a neural network to figure out what numbers should associate with each word, the first thing we do is create inputs for each unique word. In this case we have 4 unique words, so we have 4 inputs. we assign random weights to each word. (for simplicity, each word is in 2-dimensional space).



- using back-propagation, the weights of each word will be adjusted, such that similar words in similar contexts represent similar embeddings

After training

Before training



→ now, our model, knows both "troll2" and "GymRatzi" are names of movies and that they should lie within same context.

Note: that model is designed to predict next word. for example
feeding "is" into our NN, it will predict next word is "great".
However, the main goal with this model, is to adjust weights, such that similar words close together.

after the creation of this model, we will create our batch data for training. we may use "SGNS" or "CBOW"

skip-gram

- generate positive pairs
- generate k random negative pairs
- train with NN

CBOW

- instead of neighbour given center, predicts a center given its neighbours.

Date S G N S , training, updating, negative samples Subject

The core idea is:

"we want the embedding vector of word w to be more similar to the vectors of words that actually appear near it, and less similar to random unrelated words."

و د صنعت فیل، توپیں دام SGNS دیجیتی کارم کند کا امام اس پرستش -
و د چوکھا میدار و قیچی دلیم سل نہ روپڑی SGNS اوند توپیں دام (L) train
کند، چوری با فنہیں تقاویں و نام بروٹا بورن، با راستهای
سل نہ آئیں ہے منہیں گل

- ① First we see that with a sliding window of size of 2 we generate context words, how? we will create positive words pairs that are in the form (target, context1), (target, context2) with a window of size 2 that is [context1, target, context2]

- ⑦ Now we generate k random negative pairs that is a target word with all words that are not in context. ((target, and a word out of window))

important note:

We don't pick those samples uniformly. We pick them according to their unigram frequency $p(w)$ raised to power α . This makes common words less likely to be dominant. Makes rare words a bit more likely as negative.

$$p_d(w) = \frac{\text{count}(w)}{\sum_{w'} \text{count}(w')^\alpha}$$

III Date Prediction task:

Subject

input : a pair (w, c) (target, context)

label : ① if it's real positive pair, 0 if it's negative.

IV Update embedding:

inside the model, each input pair has two learnable vectors:

① a target word vector embeddings

② a context word embedding.

we compute dot product of these vectors, pass it through a sigmoid function to get probability :

$$P(\text{context is real} | w, c) = \sigma(\bar{v}_w \cdot v_c)$$

then we use loss function (binary cross entropy) to adjust:

$$L = -\lg(\sigma(v_w \cdot v_c)) - \sum_{i=1}^K \lg(\sigma(-v_w \cdot v_i))$$

this says:

- increase the dot product with real context words.

- decrease the dot product with negative words

make vectors more or less similar.

dot product as similarity signal:

if \vec{v}_w and \vec{v}_c point same direction, their dot product is large positive.

if they are near orthogonal (unrelated), dot product is near zero

if they point oppositely, dot product is negative.