

Word Representation

A perfect word representations should have:

→ Shared Lemmas:

Word with same root should have similar representation.

Ex: Mouse / mice , class / classes , eat / ate

→ Different word sences:

Should be able to distinguish between different meaning of the same word.

Ex: River Bank, Money Bank, Blood Bank
Computer mouse , pet mouse

→ Synonyms:

Synonym words should be close in the vector space.

Ex: Couch / Sofa , Car / automobile

These vectors should have high cosine similarity.

→ Antonyms:

Antonyms still be related but in opposite direction

Ex: long / short , dark / light

→ Word Similarity:

Words that are not synonyms but are similar in category or roll should be close in vector space.

Ex: Cup / Coffee , doctor / nurse

→ Word Relatedness:

Word that often occur together or are functionally connected should also be related, even if they aren't similar.

Example: Cup/coffee , scapel/surgeon

→ Word valence:

Represents how positive or negative a word feels. (Emotion Polarity)

Example: High valance: Excited , relaxed (Positive emotion)

Low valance: Angry , depressed (negative emotion)

→ Word Arousal:

Represents the intensity/energy or level of activation in emotion.

Example:

High Arousal: Excited , Angry

Low Arousal: Relaxed, Depressed

* A perfect word embedding understands grammar (lemmas), meaning (senses, synonyms), relationship (relatedness, antonyms) and emotional context (valence, arousal)

* We can guess the semantics of a word from its neighbouring word.

Term Term Matrix

A term-term co-occurrence matrix X is a $|V| \times |V|$ matrix where:

- $|V|$ is the number of words in the vocabulary
- each cell $X_{i,j}$ records how often word j occurred in the context of word i
- each row X_i is the vector representation for word i

Context may be defined in different ways:

The same document

The same sentence

Within $\pm n$ words of each other

V is typically the 10,000 - 50,000 most frequent words

Each word is represented by a large vector

Example:

$d_1 = \begin{matrix} \text{any} & \text{big} & \text{cat} \end{matrix}$

$d_2 = \begin{matrix} \text{big} & \text{cat} \end{matrix}$ ●

$d_3 = \begin{matrix} \text{cat} & \text{dog} & \text{cat} \end{matrix}$

Context ± 1 word of each other

| | any | big | cat | dog |
|-----|-----|-----|-----|-----|
| any | 0 | 1 | 0 | 0 |
| big | 1 | 0 | 2 | 0 |
| cat | 0 | 2 | 0 | 2 |
| dog | 0 | 0 | 2 | 0 |

Context ± 2 word of each other

| | any | big | cat | dog |
|-----|-----|-----|-----|-----|
| any | 0 | 1 | 1 | 0 |
| big | 1 | 0 | 2 | 0 |
| cat | 1 | 2 | 2 | 2 |
| dog | 0 | 0 | 2 | 0 |

When context window is the entire document, the easy way to build term term matrix is :

- ① Build a binary BOW for the sentences
- ② Transpose it
- ③ Multiply it with the original matrix

Example:

$$\begin{aligned} d_1 &= \text{any } \text{big } \text{cat} \\ d_2 &= \text{big } \text{cat} \\ d_3 &= \text{cat } \text{dog } \text{cat} \end{aligned}$$

① Binary BOW:

$$\begin{array}{c|ccccc} & \text{any} & \text{big} & \text{cat} & \text{dog} \\ \hline d_1 & 1 & 1 & 1 & 0 \\ d_2 & 0 & 1 & 1 & 0 \\ d_3 & 0 & 0 & 1 & 1 \end{array}$$

$$\therefore M = \begin{vmatrix} 1 & 1 & 1 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & 0 & 1 & 1 \end{vmatrix}$$

$$M^T = \begin{vmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \\ 0 & 0 & 1 \end{vmatrix}$$

$$M^T \cdot M = \begin{vmatrix} 1 & 1 & 1 & 0 \\ 1 & 2 & 2 & 0 \\ 1 & 2 & 3 & 1 \\ 0 & 0 & 1 & 1 \end{vmatrix}$$

$$\begin{array}{cccc} & \text{any} & \text{big} & \text{cat} & \text{dog} \\ \text{any} & 1 & 1 & 1 & 0 \\ \text{big} & 1 & 2 & 2 & 0 \\ \text{cat} & 1 & 2 & 3 & 1 \\ \text{dog} & 0 & 0 & 1 & 1 \end{array}$$

Here,
Diagonal values represents
how many document had
the word.

→ dog appeared in 1 document.

Comparing Word vectors:

Most common similarity measure is cosine similarity.

$$\cosine(v, w) = \frac{v^T w}{\|v\| \|w\|} = \frac{\sum_i v_i w_i}{\sqrt{\sum_i v_i^2} \cdot \sqrt{\sum_i w_i^2}}$$

range of cosine similarity [-1 to 1]

$$u = [0 \ 1 \ 0 \ 1]$$

$$v = [1 \ 0 \ 1 \ 0]$$

$$w = [3 \ 0 \ 3 \ 0]$$

$$\cosine(u, v) = \frac{0*1 + 1*0 + 0*1 + 1*0}{\sqrt{0^2 + 1^2 + 0^2 + 1^2}} \cdot \frac{1^2 + 0^2 + 1^2 + 0^2}{\sqrt{1^2 + 0^2 + 1^2 + 0^2}}$$

Cosine value 1 or closer to 1 means words are similar

Cosine " 0 " " 0 " " are not related

Cosine " -1 " " -1 " " are opposite.

Term Term matrix issues:

- Very sparse
- Doesn't carry any contextual information
- Doesn't represent how a word is in a sentence.

How do we use word vectors?

Word level task: finding synonyms via cosine

classify when input is one word

For sentence & document tasks:

- Combine all word vectors by either concatenating or using centroid technique
- Or, Using RNN (More will be discussed later.)
- Then these vectors can be used for classification.

Sparse vs Dense vector:

So far, all of the techniques we learnt were sparse.

Sparse vector & matrices have their own disadvantages. So, we will now focus on dense vectors:

Dense vector has following advantages:

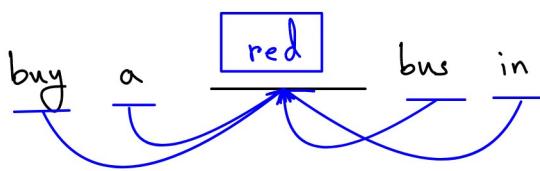
- fewer weights to learn
- fewer features can reduce overfitting
- forces sharing as there are not enough dimensions for each word to be completely independent.

Word Embedding:

→ Instead of looking at co-occurrence, let's try to predict

Two ways:

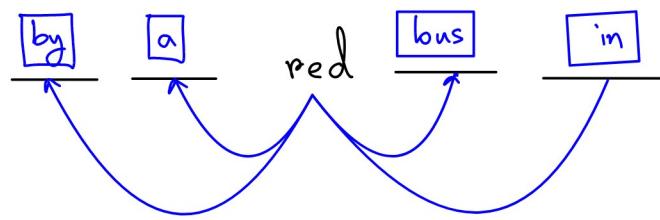
CBOW: Given neighbour words, predict the target word



$$\text{context/window_size} = 2$$

means how many neighbour
we will consider

Skipgram: Predict neighbouring words given the target word.



Cbow is easy, lets focus on Skipgram.

Creating training data is very easy. Given a corpus, we select each word and its neighbour words (each a pair) and append them in training data.

Example: Content window = 2

g live in Dhaka....

| input word | target word |
|------------|-------------|
| g | live |
| g | in |
| live | g |
| live | in |
| live | Dhaka |
| in | g |
| in | live |
| in | Dhaka |
| in | . |
| Dhaka | live |
| Dhaka | in |
| : | : |

* For the very starting word/s and ending word/s , we can add dummy words.

Example:

| input word | target word |
|------------|-------------|
| g | <dummy> |
| g | <dummy> |

Initial idea of training an embedding model.

- ① for each word in the dataset, create a one hot vector (shape $1 \times |V|$)
- ② Pass it to the neural network to calculate prediction
- ③ Project to output vocabulary using softmax activation function.

The NN architecture:

Vocabulary: { f , live, in, Dhaka }

f live in Dhaka

$$\text{f} = [1 \ 0 \ 0 \ 0]$$

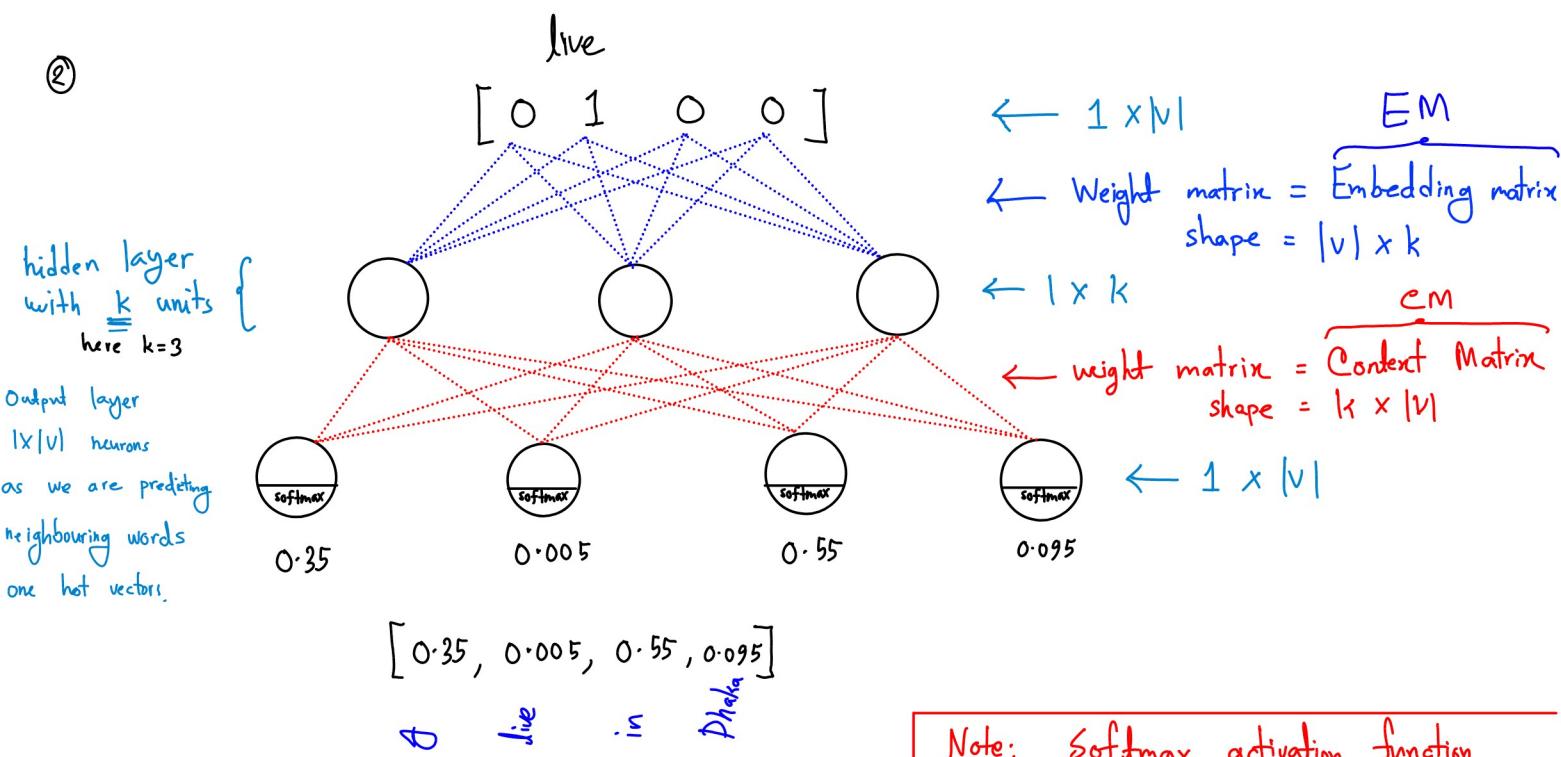
① One hot vector for each word is = live = $[0 \ 1 \ 0 \ 0]$

$$\text{in} = [0 \ 0 \ 1 \ 0]$$

Let, Our training data is (live, in)

$$\text{Dhaka} = [0 \ 0 \ 0 \ 1]$$

②



③

$$\text{prediction of NN} = [0.35, 0.005, 0.55, 0.095]$$

$$\text{Actual target (in)} = [0, 0, 1, 0]$$

$$\text{Error} = \text{Actual} - \text{target}$$

$$= [-0.35, -0.005, 0.45, -0.095]$$

Now from this error update weight matrices.

Note: Softmax activation function

converts a tuple of K real numbers into a probability distribution of K outcomes.

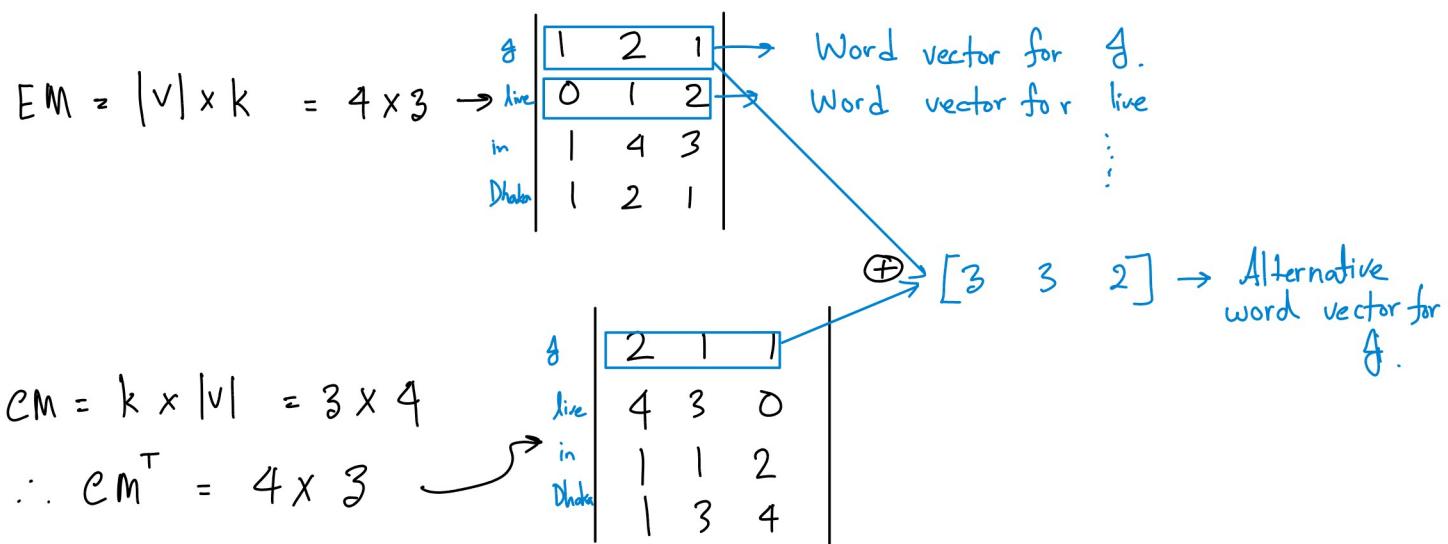
Output Layer:

$$y = \begin{bmatrix} 1.3 \\ 5.1 \\ 2.2 \\ 0.7 \\ 1.1 \end{bmatrix} \rightarrow \text{Softmax} \rightarrow \begin{bmatrix} 0.02 \\ 0.90 \\ 0.05 \\ 0.01 \\ 0.02 \end{bmatrix}$$

$$\text{softmax}(y_i) = \frac{e^{y_i}}{\sum_{j=1}^n e^{y_j}}$$

$$\text{sum} = 1$$

After training this Neural Network on our entire dataset, we will get optimum weight matrices: Embedding & Context Matrices. Let, for the above network, we got,



Note that: We can either use only EM for our word vector or use the summation of EM & CM for our word vectors.

Also, we are getting a $1 \times k$ dimension (1×3 in this example) vectors.

So, how many neuron(k) we will keep in our hidden layer will determine the dimension of the vectors we get. Usually we keep $50 \leq k \leq 1000$.

Theoretically, this model should give us properly learnt word vectors.

But there is a problem. The step (3) ↴

③ Project to output vocabulary using softmax activation function. is extremely computationally expensive. In proper dataset & training environment, $|V|$ can be from 50K to 1 million. This means the output layer has too many neurons & for each of them we need to compute softmax "for each training data."

This is where Computer Scientist of Czech, Mikolov entered in the game.

Mikolov reformulated the task. Instead of predicting neighbouring word, Mikolov converted it to a simple classification task.

Before:

f (live, in) Dhaka
After:



that means, given two words from the training data, we try to predict whether they are neighbour or not.

Steps:

- ① Converting the training data

Example: Context window = 2

f live in Dhaka....

[Word2Vec Training Steps]

| input word | Output word | target |
|------------|-------------|--------|
| f | live | 1 |
| f | in | 1 |
| live | f | 1 |
| live | in | 1 |
| live | Dhaka | 1 |
| : | : | : |
| f | Dhaka | 0 |
| f | f | 0 |
| Dhaka | f | 0 |
| f | else | 0 |

- ② Negative Sampling:

After step 1, all of the target in training data is 1. All 1 in training data means all words in training data are neighbour which is not good for learning. So, we randomly add some words which are not neighbour with target = 0.

This is called random sampling.

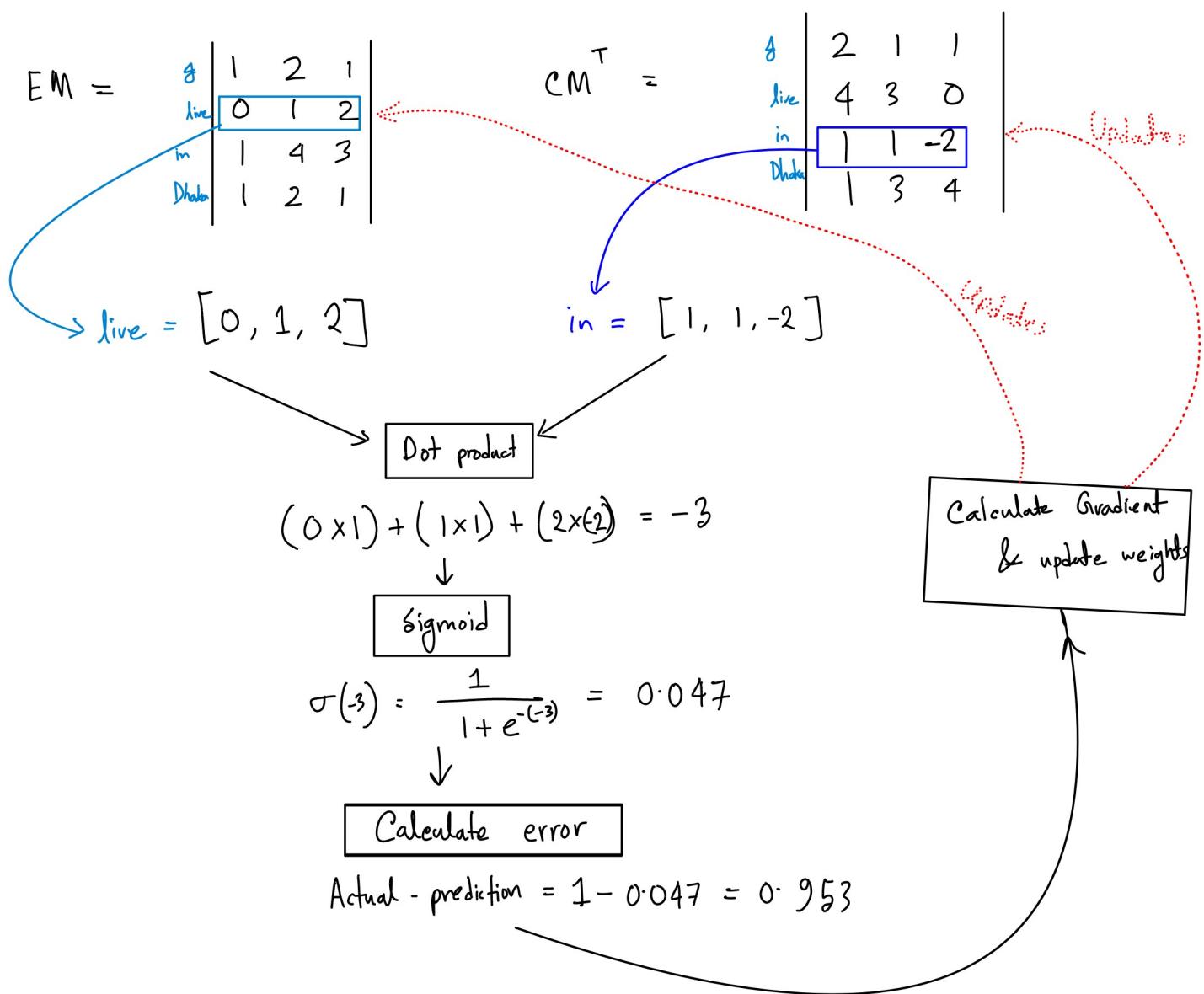
- ③ Initialize Embedding Matrix ($|V| \times k$) & Context Matrix ($k \times |V|$) with random values.

- ④ Take input word vector from EM and output word vector from CM
- ⑤ Dot product two vectors and pass them on sigmoid function for binary classification.
- ⑥ Calculate error and update the word vectors accordingly on both EM & CM.

Doing this for the entire dataset will give us the trained EM and CM.

Example: Let,

Training data = (live, in, 1)



This entire process gives us the updated weight Matrix EM & CM.
Now we can use them as word vectors.

While working, we do not need to code from scratch. We can directly use gensim's Word2vec.

We don't need to train Word2vec always. It is rarely necessary to train skipgram. Instead, we trained it once on massive data and use this pre-trained embeddings directly.

Semantic Properties of Embeddings

- * Different type of similarity or association

Based on the context window, word association changes:

Smaller window size (2 - 15) embeddings

→ high similarity score indicates that words are interchangeable

Example: Hogwart is similar to Sunnydale

(both are school/
academy)

→ Synonym or similar word stays close in the vector space.

Larger window size (15-50 or more) embeddings

→ high similarity score indicates relatedness.

Example: Hogwart is more similar to dumbledore, malfoy etc.

(As they are from
same book/stories.)

- * Analogy / Relational Similarity:

Embeddings captures relational meanings. Therefore, we can do analogies like:

$$\begin{array}{ccccccc} \xrightarrow{\quad} & \xrightarrow{\quad} & \xrightarrow{\quad} & & \xrightarrow{\quad} \\ \text{King} & - & \text{man} & + & \text{woman} & \approx & \text{Queen} \\ \xrightarrow{\quad} & \xrightarrow{\quad} & \xrightarrow{\quad} & & \xrightarrow{\quad} \\ \text{Paris} & - & \text{France} & + & \text{Italy} & \approx & \text{Rome} \end{array}$$

Historical Semantic Context:

The data we use reflects on our embedding.
So, we can see how the meaning of word changes through time.

We can visualize this using PCA or t-SNE which compresses the high dimensional word vectors into lower dimension for plotting. (Ex: 300 D to 2D)

Bias in Embeddings:

As embeddings are learnt from training data, it captures the bias from it as well.

Example: $\begin{cases} \text{cosine}(\text{man}, \text{attractive}) = 0.3085 \\ \text{cosine}(\text{woman}, \text{attractive}) = 0.4111 \end{cases}$
 $\begin{cases} \text{cosine}(\text{dumb}, \text{American}) = 0.4118 \\ \text{cosine}(\text{dumb}, \text{European}) = 0.2658 \end{cases}$

Also,

he is a nurse . she is a farmer.



ଯେ ଏକଜନ ମାର୍ଗୀ , ଯେ ଦ୍ୱିତୀୟ କୁଷକ !



↓

Translation Model

$B \rightarrow E$

↓
She is a nurse. He is a farmer.

Static Word Embedding:

So far, the word vectors we learnt about has a fixed vector after training. They ignore the context in that sentence.

Ex: The river bank, A bank deposite

$$\begin{bmatrix} 1 \\ 2 \\ 1.3 \\ \vdots \end{bmatrix}$$

$$\begin{bmatrix} 1 \\ 2 \\ 1.3 \\ \vdots \end{bmatrix}$$



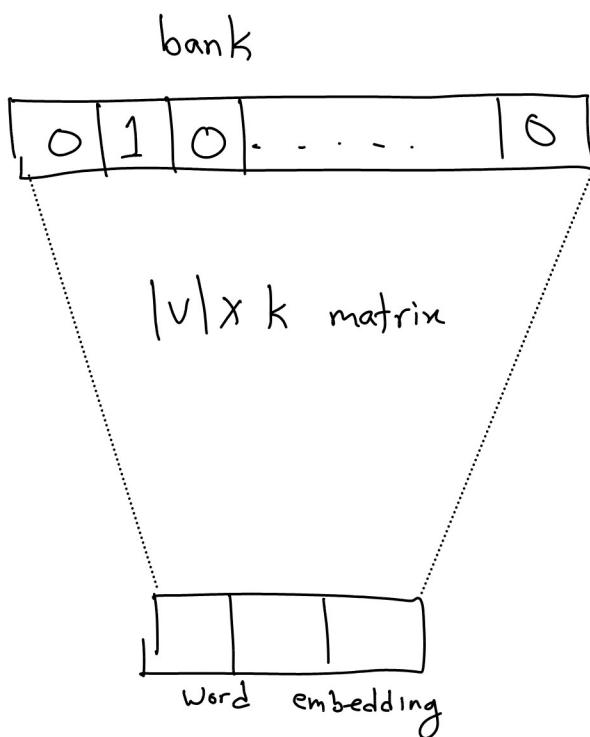
They shouldn't be similar,
Based on context, they
should be different.

Static Embedding captures both the meaning of money bank and river bank. But this is not accurate.

Solution ?

Contextual Word Embeddings:

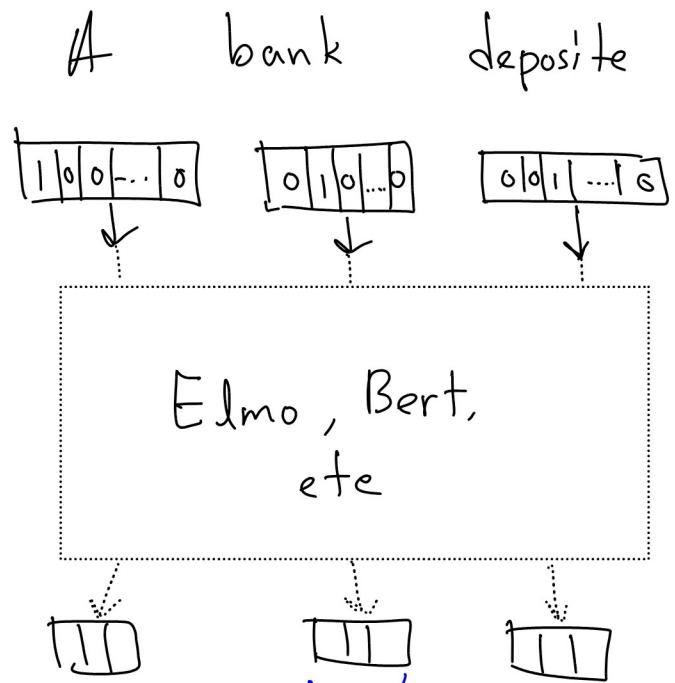
Static Word Embedding:



Input 1 word

Output 1 embedding

Contextual Word Embedding:



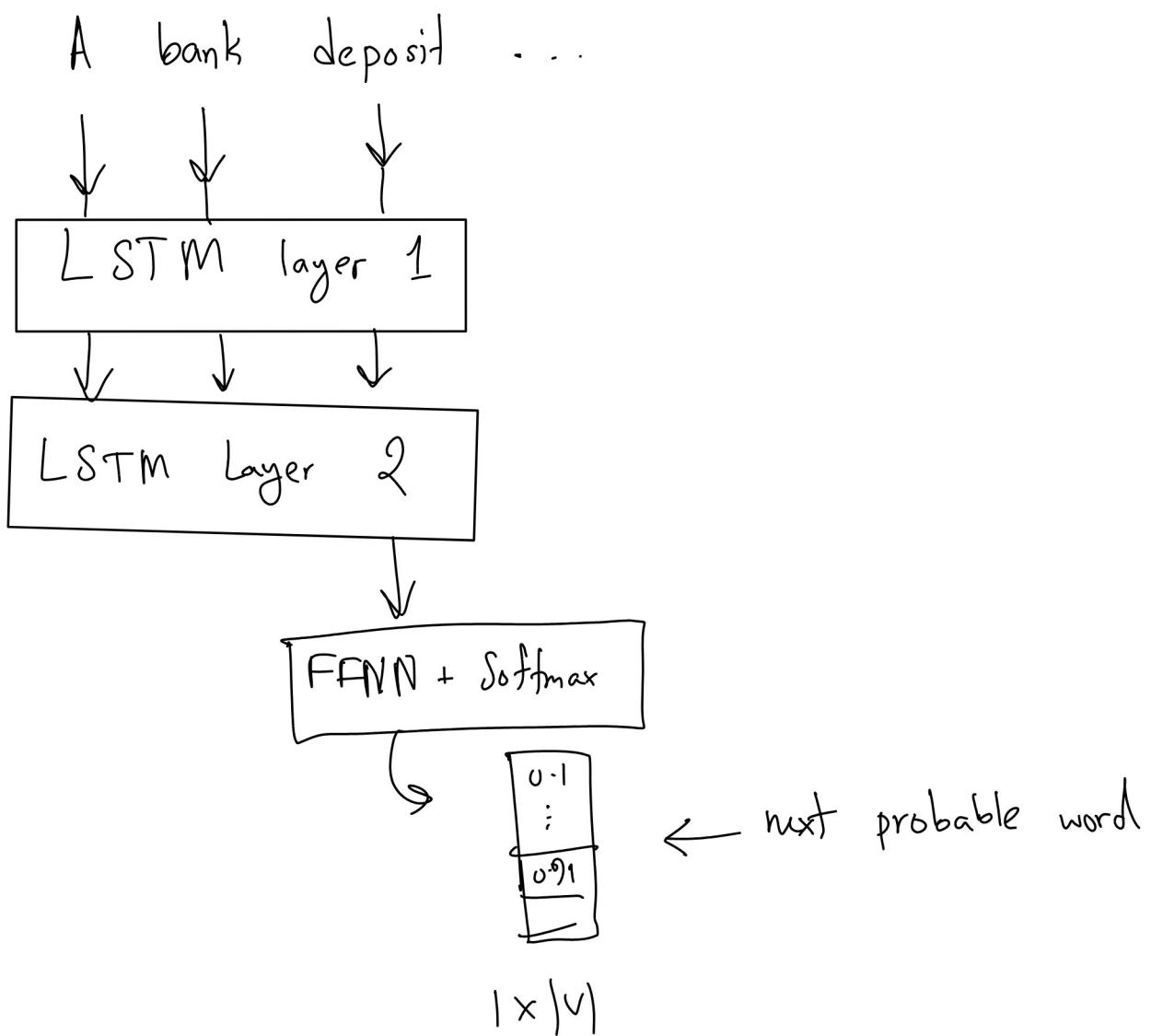
Input N word
Output N embedding

→ This 'bank' embedding has the context of neighbour words. So, it has the properties of money bank.

Language Model:

- Given a sequence of word, predict the next word.
- Unsupervised, great for learning word representation.

Elmo's task:



How to use contextual embedding?

- ① Train on Unlabeled data } Pre training
- ② Extract word vectors to use as feature
- ③ Fine tune it on labeled data. } finetuning,