# How do you Learn Embedding

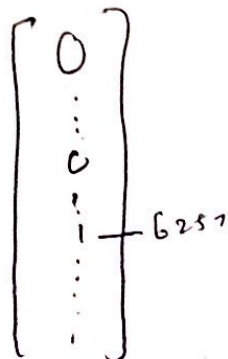## Embedding Matrix.

Let's say we are working on a 10,000 word vocabulary.

a   aaron ..... Orange ........ Zulu   ← 6257



$300$

$\longleftarrow$ 6257

1st row
2nd row

$E \rightarrow$ as Embedding Matrix

$O_{6257} \longrightarrow$ 1-hot vector

10,000

$c_w =$ Embedding for word$_w$

$E \cdot O_{6257} =$

$300 \times 10,000$  dot product  $\times 1$

$\longrightarrow$ cloumn in Orange Column
$\longrightarrow$ 2nd row

$(300, 1)$

$= e_{6257}$  $300 \times 1$ Vector

$\longrightarrow$ 300-dimension Vector correspanly to word orange

More generally $e$ times $O$

$$E_j \cdot O_j = e_j \quad \text{embedding for word}_j$$

goal is to learn $E_j$

1- Initialize randomly and learn all
   parameters of this 300 by 10000 dimen
   matrix.

... A Language model (LM) is the
   use of various statistical and probablistic
   techniques to determine the probability
   of a given sequence of words.
   occuring in a sentence.

   Lets assume we are building
   an LM. using NN

Task:.

   I want a glass of Orange ___

   4343 9665  I   3852  6163  6257   .
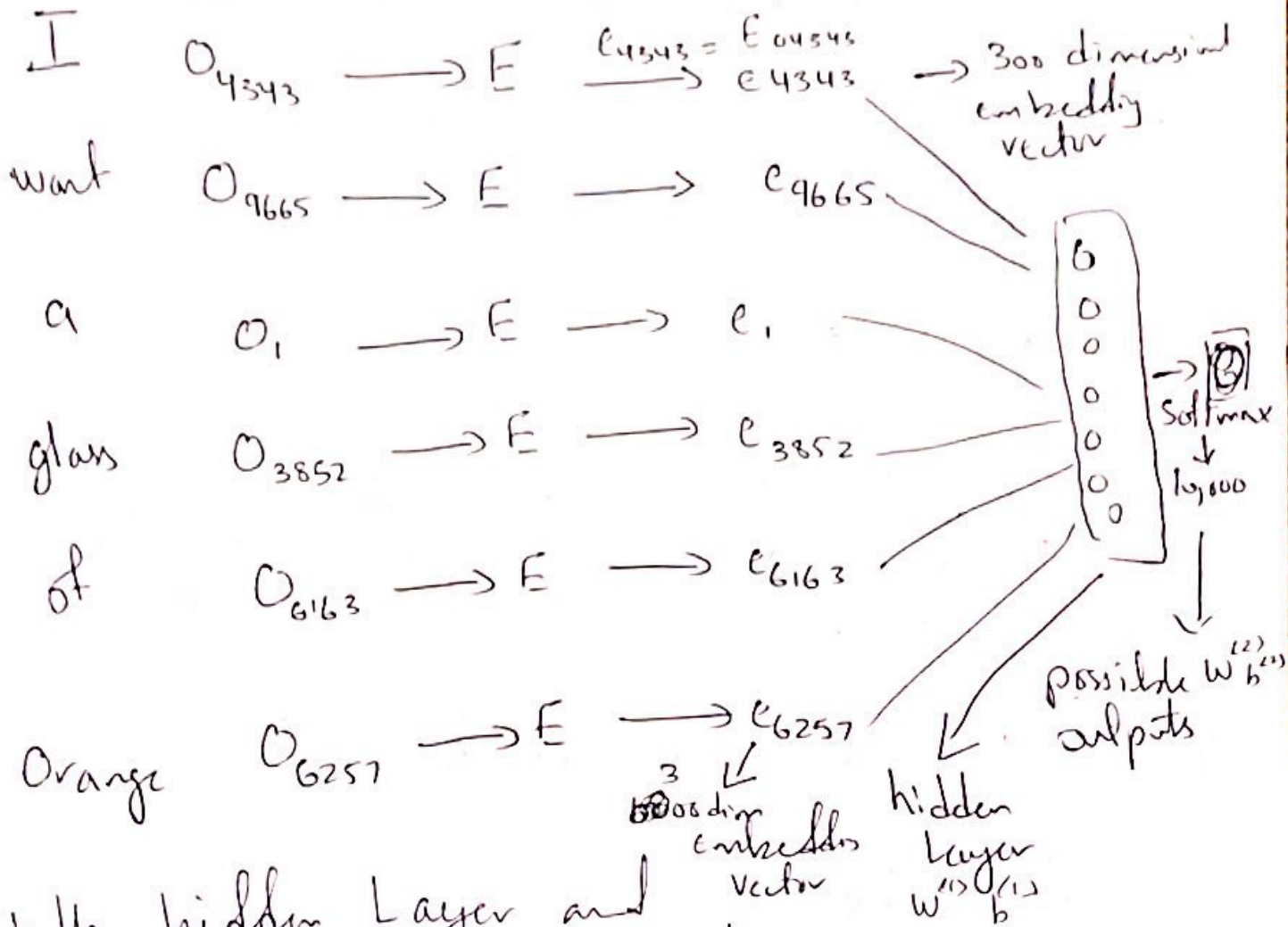
   predict the next word          index in
                                   vocabulary.

Predicting the next word in the sequence is a reasonable way of learning a set of embeddings.

I want a glass of Orange _____
4343 9665 1 3852 6163 6257

I $O_{4343}$ ⟶ E $\xrightarrow{e_{4343} = E \cdot o_{4545}}$ $e_{4343}$ ⟶ 300 dimensional embedding vector

want $O_{9665}$ ⟶ E ⟶ $e_{9665}$

a $O_1$ ⟶ E ⟶ $e_1$

glass $O_{3852}$ ⟶ E ⟶ $e_{3852}$

of $O_{6163}$ ⟶ E ⟶ $e_{6163}$

Orange $O_{6257}$ ⟶ E ⟶ $e_{6257}$

6
0
0
0
0
0
0

⟶ Softmax ↓ 10,000

possible $W^{(2)} b^{(2)}$ outputs

300 dim embedding vector
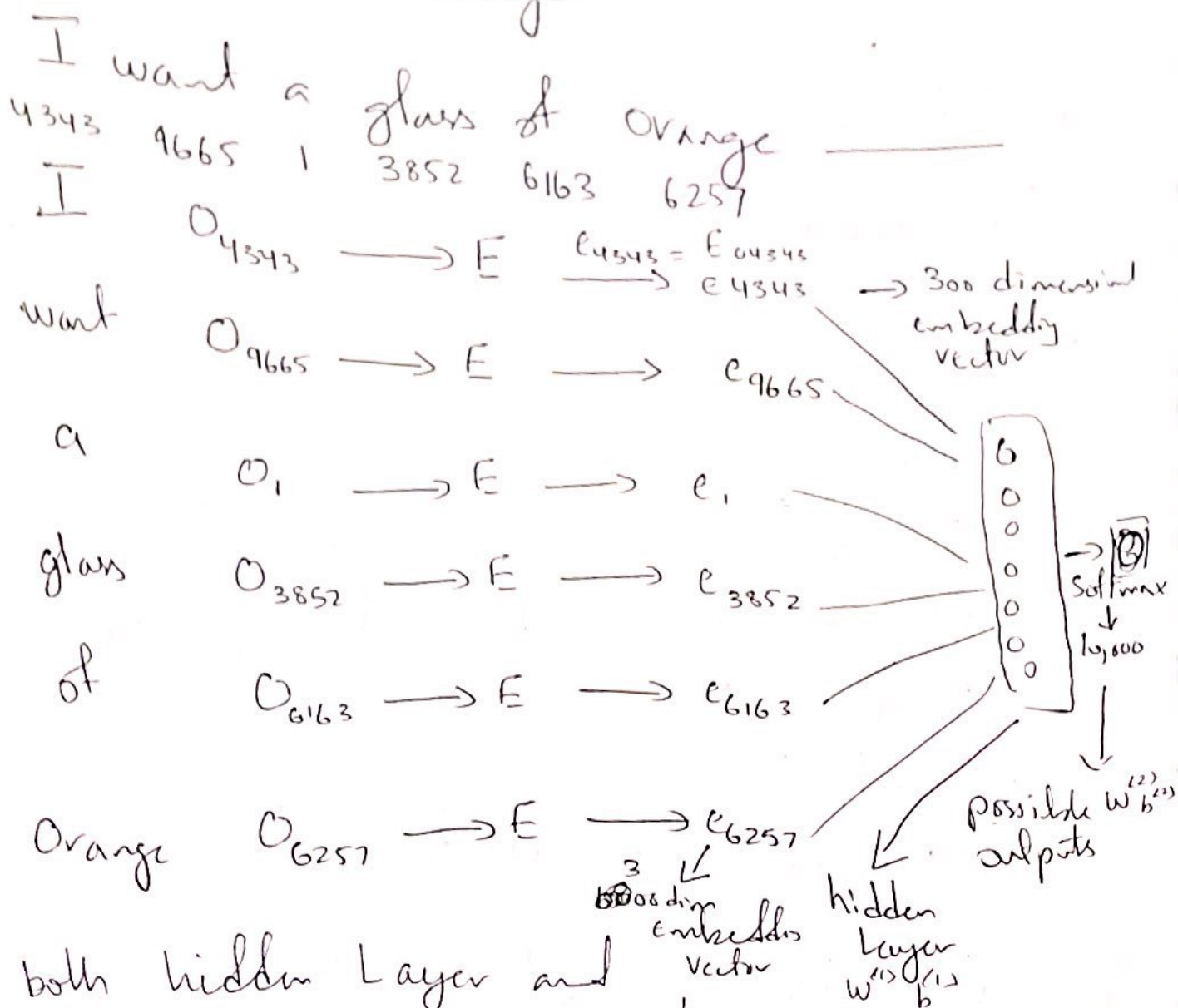
hidden layer $W^{(1)} b^{(1)}$

both hidden Layer and the softmax layer will have their own params.

1800 input layer
6 x 300 stacking them together.

Predicting the next word in the sequence is a reasonable way of learning a set of embeddings.

I want a glass of orange ___
4343  9665  1  3852  6163  6257

$I \quad O_{4343} \longrightarrow E \xrightarrow{\ e_{4343} = E \cdot O_{4343}\ } e_{4343} \longrightarrow$ 300 dimensional embedding vector

want $\quad O_{9665} \longrightarrow E \longrightarrow e_{9665}$

a $\quad O_1 \longrightarrow E \longrightarrow e_1$

glass $\quad O_{3852} \longrightarrow E \longrightarrow e_{3852}$

of $\quad O_{6163} \longrightarrow E \longrightarrow e_{6163}$

Orange $\quad O_{6257} \longrightarrow E \longrightarrow e_{6257}$

$\rightarrow$ Softmax $\downarrow$ 10,000

possible outputs $W^{(2)} b^{(2)}$

hidden Layer $W^{(1)} b^{(1)}$

3 $\swarrow$ 1800 dim embedding vector

both hidden Layer and the softmax layer will have their own params.

1800 input layer
6 x 300 stacking them together.

have a fixed historical window

## a glass of orange ——

predict the next word given previous
four words.

- In such a case NN will input
- $4 \times 300 = 1200$ dimensional input feature
  vector

- Using a fixed window means
  we can deal with arbitrary long
  sentences because input size is
  fixed.

- Us parameter matrix $'E'$, same
  matrix for all the words.

  $w^{\langle 1 \rangle}, b^{\langle 1 \rangle}$ ad $w^{\langle 2 \rangle} b^{\langle 2 \rangle}$

  $w, b$ are also parameters of the
  algorithm.

- Use back propagation to perform
  gradient descent to predict,

given 4 words in the sequence what will be the next word in the text.

This algorithm will perform reasonably well.

Example: Orange juice

Apple juice

Algo will give similar weights to orange and Apple.

if we have 300 features the algo will ~~fit better w~~ find it best if fruits will end up with similar feature vectors

# Other Algorithm Complex Example

I want a glass of orange juice to go
along with my cereal.

a) Context: Last 4 words

b) 4 words left and right

a glass of orange ___?___ to go along
with

c) Last 1 word

orange ___?___

NN with just 1 word.

d) take nearby 1 word

glass ___?___

this is an idea of skip gram

for Language model: use the last
few words. as context
to learn word embedding. use all of
these contexts.

Two ways to implement word2vec, CBOW,
and skip-gram

CBOW: a window around some target
word and consider the words around it.
(context)

We supply context words and
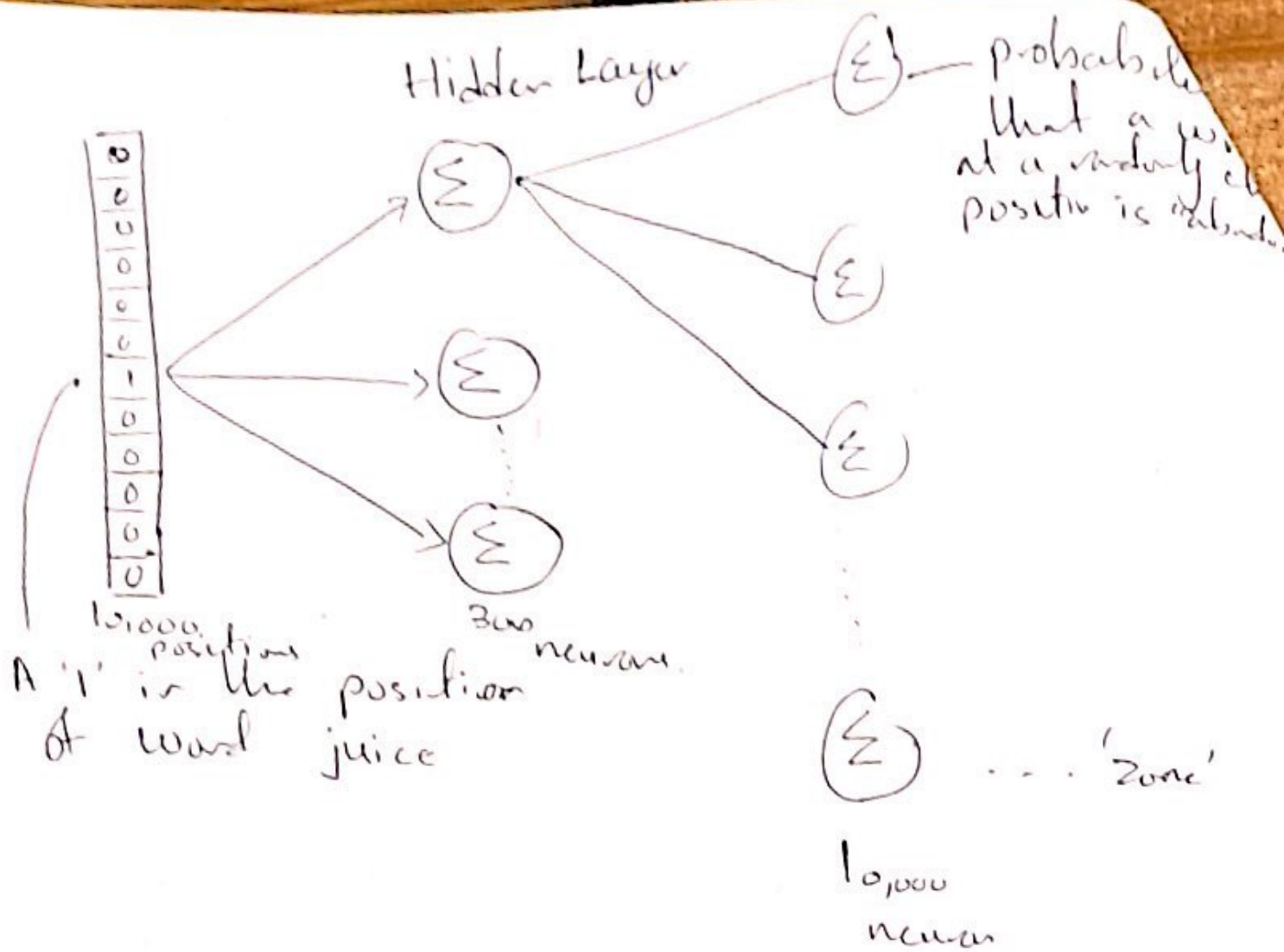use it to predict the target word.

skip-gram: We have a target word
and we try to predict the words that
are in the window around that word.
i.e. predict the context

The input words are passed in
as one-hot encoded vectors.
Then we will send it to a
hidden layer then to a softmax layer
The weights of the hidden layer
represents word embeddings

Hidden Layer



— probability that a word at a randomly ch position is [unreadable]

A '1' in the position of word juice

10,000 positions

300 neurons

'zone'

10,000 neurons

The embedding matrix has a size of number of words by the number of neurons in the hidden layer. So for 10,000 words and 300 hidden units, the matrix will have 10,000 x 300

$$\begin{bmatrix} 0 & 0 & 0 & 1 & 0 \end{bmatrix} \times \begin{pmatrix} 17 & 24 & 1 \\ 23 & 5 & 7 \\ 4 & 6 & 13 \\ 10 & 12 & 19 \\ 11 & 18 & 25 \end{pmatrix} = \begin{bmatrix} 10 & 12 & 19 \end{bmatrix}$$

for the fourth word

# Word 2 vec Model

skip-gram

I want a glass of orage juice to go along with my cereal.

| Context | Target |
|---|---|
| Orange | |