

NLP Basics:-

Count Based Methods → Count vectorizer → One hot encoding vectors

Corpus:- (1) The food is awesome → ['food', 'awesome']

(2) Ambience is good, food is good → ['ambience', 'good', 'food']

(3) Ambience is good, but food is bad → ['ambience', 'good', 'food', 'bad']

Vocabulary = $V = \{ \text{'food', 'awesome', 'ambience', 'good', 'bad'} \}$

One hot Encoding

	'food'	'awesome'	'ambience'	'good'	'bad'
①	1	1	0	0	0
②	1	0	1	1	0
③	1	0	1	1	1

Count Vectorizer

Vectors are also of the size of- vocabulary

$$C_i^{(d)} = \text{Count}(w_i, d) \quad w_i \in V$$

$$\vec{d} = [C_1^{(d)}, C_2^{(d)}, C_3^{(d)}, \dots, C_n^{(d)}]$$

With previous example the count vector will look like :-

	food	awesome	ambience	good	bad
①	1	1	0	0	0
②	1	0	1	②	0
③	1	0	1	1	1

tf-idf vectorizer:-

tf \rightarrow term frequency , idf \rightarrow inverse document frequency

document freq:- $w_i \rightarrow i^{\text{th}}$ term
 $n_i \rightarrow$ number of documents it has occurred.

$df_i = \frac{n_i}{n}$, $n \rightarrow$ total no. of documents.

$$idf_i = \log\left(\frac{n}{n_i}\right)$$

$$tf-idf(w_i, d_j) = tf(w_i, d_j) \times \log\left(\frac{n}{n_i}\right)$$

term freq:- $tf(w_i, d_j) = \text{Count}(w_i, d_j)$

Word2Vec

In word2vec the words are represented by dense vectors of similar dimension ($d=300$)

↳ CBOW → Continuous Bag of Words

↳ Skip-gram method.

CBOW

(1) Obtain one-hot encoded vector of each word in the vocabulary.

$V = \{ \text{'food', 'good', 'ambience', 'awesome', 'bad', 'wine'} \}$

$$\text{OHE}(w_i) = \begin{bmatrix} 0 & 0 & 0 & \dots & 1 & 0 & 0 \dots \dots \dots 0 \end{bmatrix}_{1 \times |V|}^T$$

$$\text{OHE}(\text{food}) = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}^T, \quad \text{OHE}(\text{ambience}) = \begin{bmatrix} 0 & 0 & 1 & 0 & 0 & 0 \end{bmatrix}^T$$

$$\text{OHE}(\text{bad}) = \begin{bmatrix} 0 & 0 & 0 & 0 & 1 & 0 \end{bmatrix}$$

(2) Obtain Context Matrix with a specified content window (k)
 $C \rightarrow$ context matrix, $C_{|V| \times |V|}$, $C_{ij} \in C \rightarrow$ How many times word w_i has appeared in the content of w_j

$k=2$

Corpus: ["I like NLP, but I hate python",
 "I like python, but I love cv"]

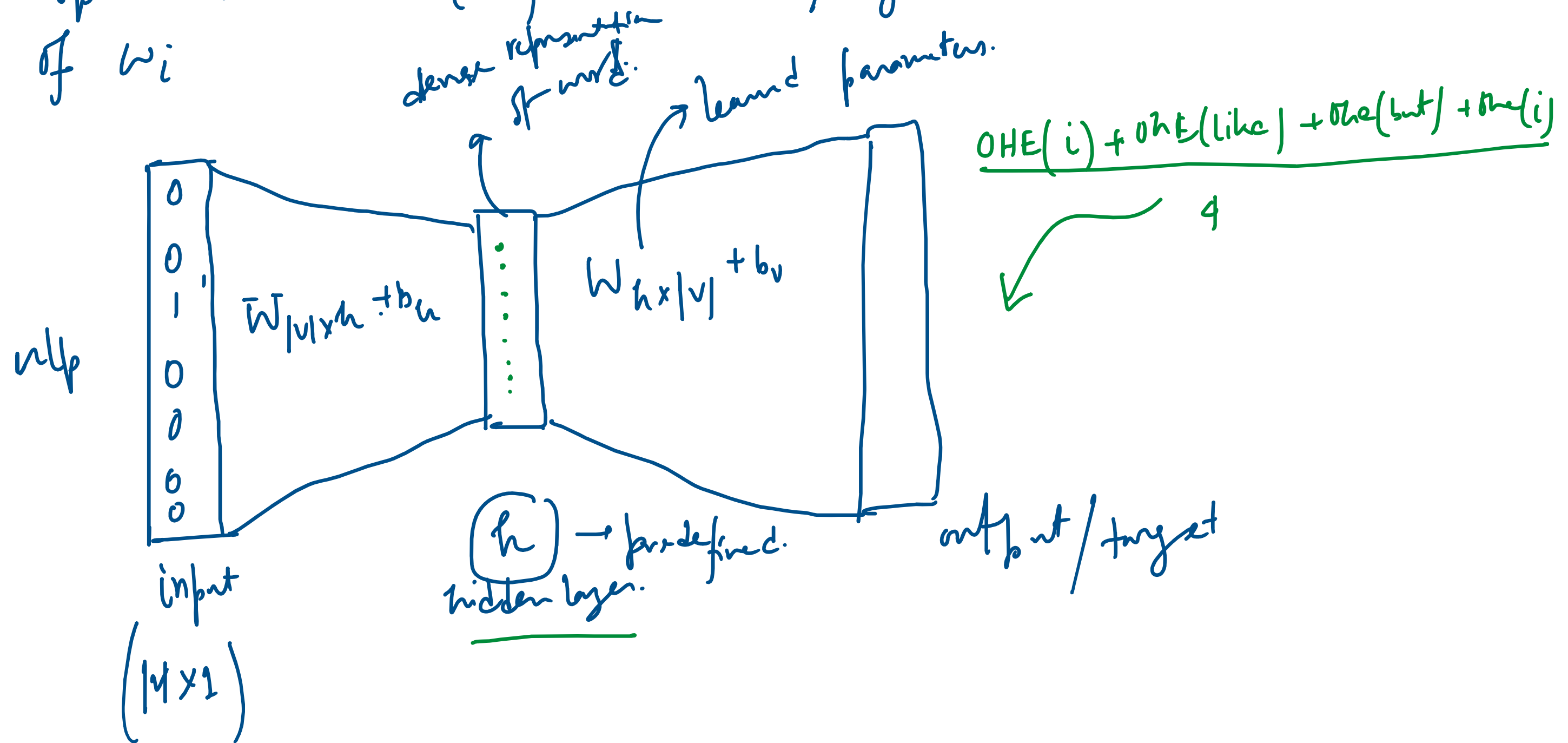
$V = \{ 'I', 'like', 'nlp', 'hate', 'python', 'love', 'cv' \}$

	'i'	'like'	'nlp'	'hate'	'python'	'love'	'cv'
i	0	2	1	1	2	1	1
like	2	0	1	0	1	0	0
nlp							
hate							
python	2	1	0	1	0	0	0
love							
cv							

$k=2$

(3) Train a shallow neural net

input is the $OHE(w_i)$ & output/target average of the context of w_i



Advantages of dense representation of word vectors over sparse representation

- 1) It takes less space to store dense representation.
- 2) It preserves semantic meaning.

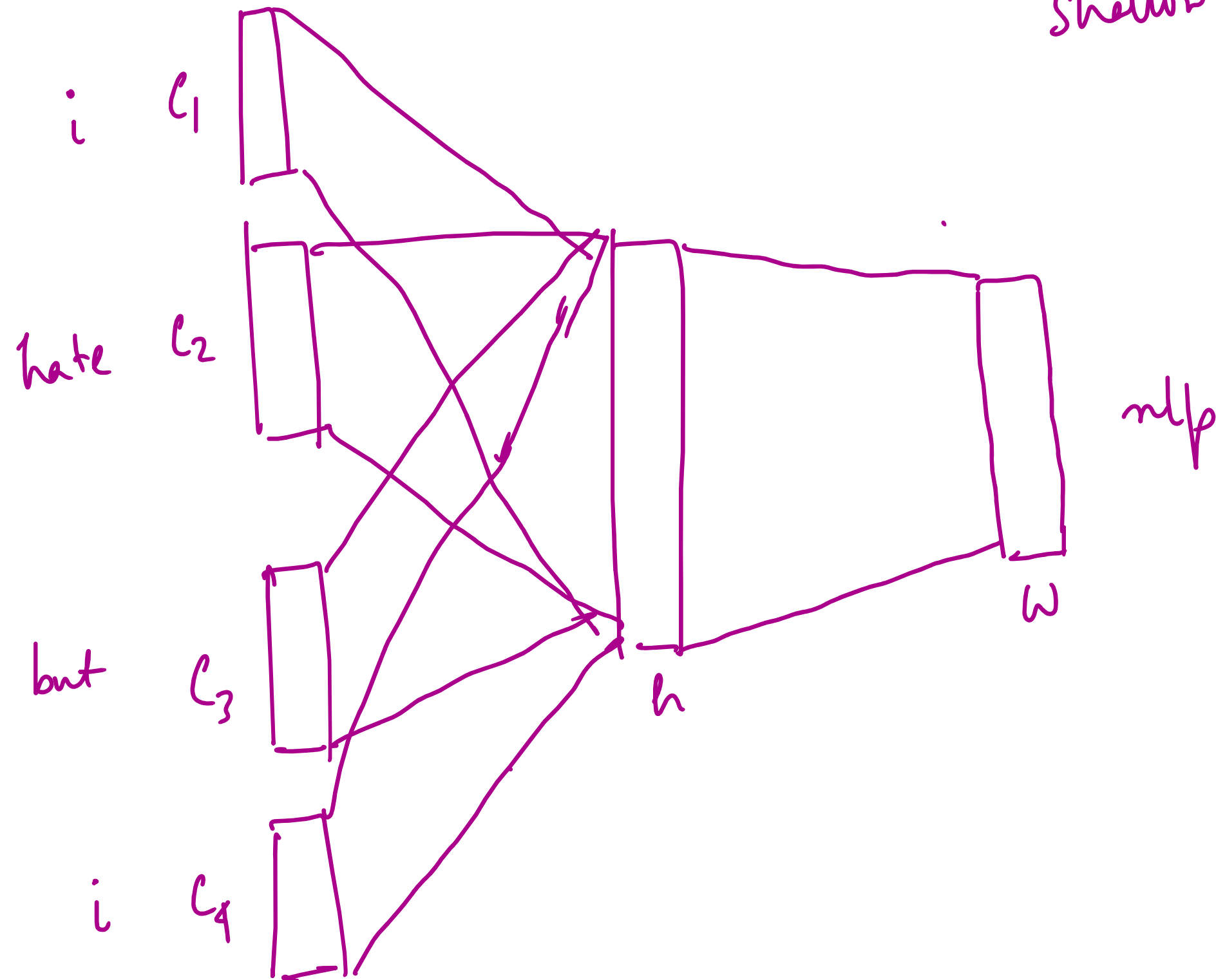
USA → washington DC

India → ? (Delhi)

$$V(\text{USA}) - V(\text{washington DC}) \approx V(\text{India}) - V(\text{delhi})$$

Skip-gram model :- $k=2$

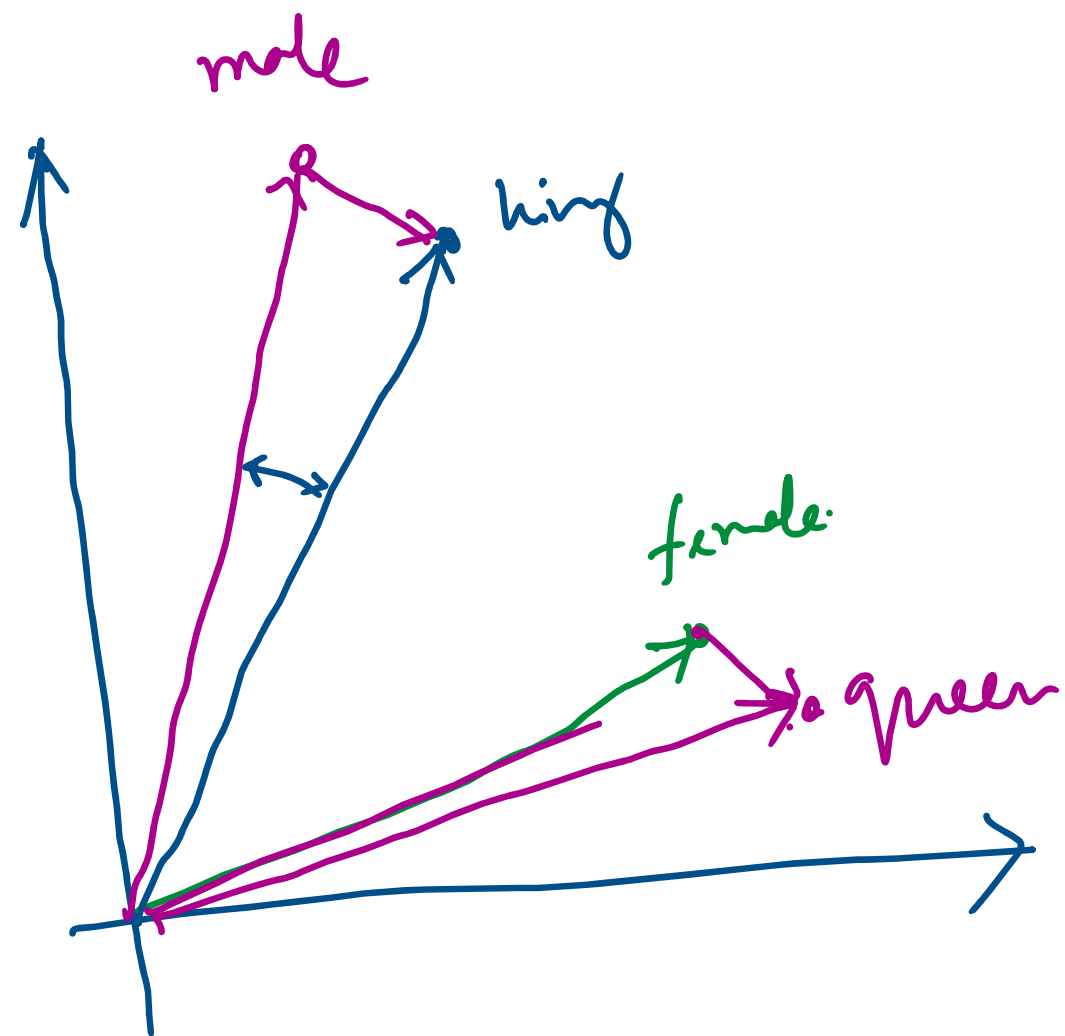
Shallow dense neural network.



Pretrained Word2Vec

- Google has trained Word2Vec models with wikipedia data.
- Using that model it has generated word-vectors for all words in English.
- It has also word-vector repository for other languages.
- These pretrained word vectors can be used directly for other NLP tasks, such as: Chatbot, Sentiment analysis etc.
- These pretrained models are available to download for free.

We can also create our own domain aware word2vec using python library 'gensim'.



king, queen.
 $\cosine\ sim(king, male)$ is higher
 compared to $\cosine\ sim(male, female)$

$$\vec{male} - \vec{king} \approx \vec{female} - \vec{queen}$$

male \rightarrow female
 king \rightarrow queen.

deer, bison, chair

