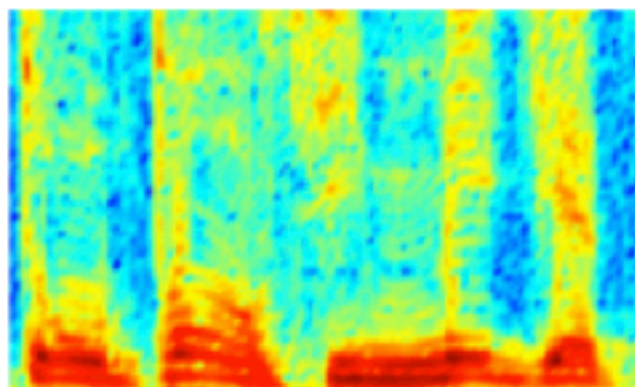


Embeddings

Subash Gandyer
Data Scientist, HealthChain

Why embeddings

AUDIO



Audio Spectrogram

DENSE

IMAGES

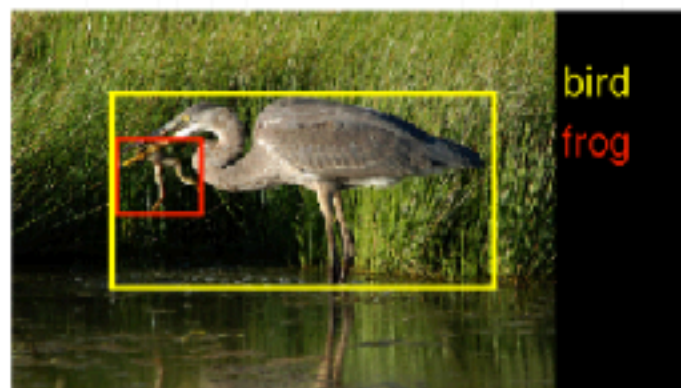


Image pixels

DENSE

TEXT

0	0	0	0.2	0	0.7	0	0	0
---	---	---	-----	---	-----	---	---	---	-----	-----

Word, context, or
document vectors

SPARSE

Word Embedding

‘the collective name for a set of language modeling and feature learning techniques in natural language processing (NLP) where words or phrases from the vocabulary are mapped to vectors of real numbers’

the vector, which reflects the structure of the word in terms of morphology (Enriching Word Vectors with Subword Information) / word-context(s) representation (word2vec Parameter Learning Explained) / global corpus statistics (GloVe: Global Vectors for Word Representation) / words hierarchy in terms of WordNet terminology (Poincaré Embeddings for Learning Hierarchical Representations) / relationship between a set of documents and the terms they contain (Latent semantic indexing) / etc

- Idea: Capture morphological / semantic / context / hierarchical information from words
- Problem: Choosing which embedding works best

Embeddings

- One-hot encoding
- TF-IDF encoding
- Word2Vec
- Glove

One-hot encoding / CountVectorizer

- Idea: Collect a set of documents (words, sentences, paragraphs, articles) Count every occurrence of word

```
from sklearn.feature_extraction.text import CountVectorizer
# create CountVectorizer object
vectorizer = CountVectorizer()
corpus = [
    'Text of first document.',
    'Text of the second document made longer.',
    'Number three.',
    'This is number four.',
]
# learn the vocabulary and store CountVectorizer sparse matrix in X
X = vectorizer.fit_transform(corpus)
# columns of X correspond to the result of this method
vectorizer.get_feature_names() == (
    ['document', 'first', 'four', 'is', 'longer',
     'made', 'number', 'of', 'second', 'text',
     'the', 'this', 'three'])
# retrieving the matrix in the numpy form
X.toarray()
# transforming a new document according to learn vocabulary
vectorizer.transform(['A new document.']).toarray()
```

TF-IDF encoding

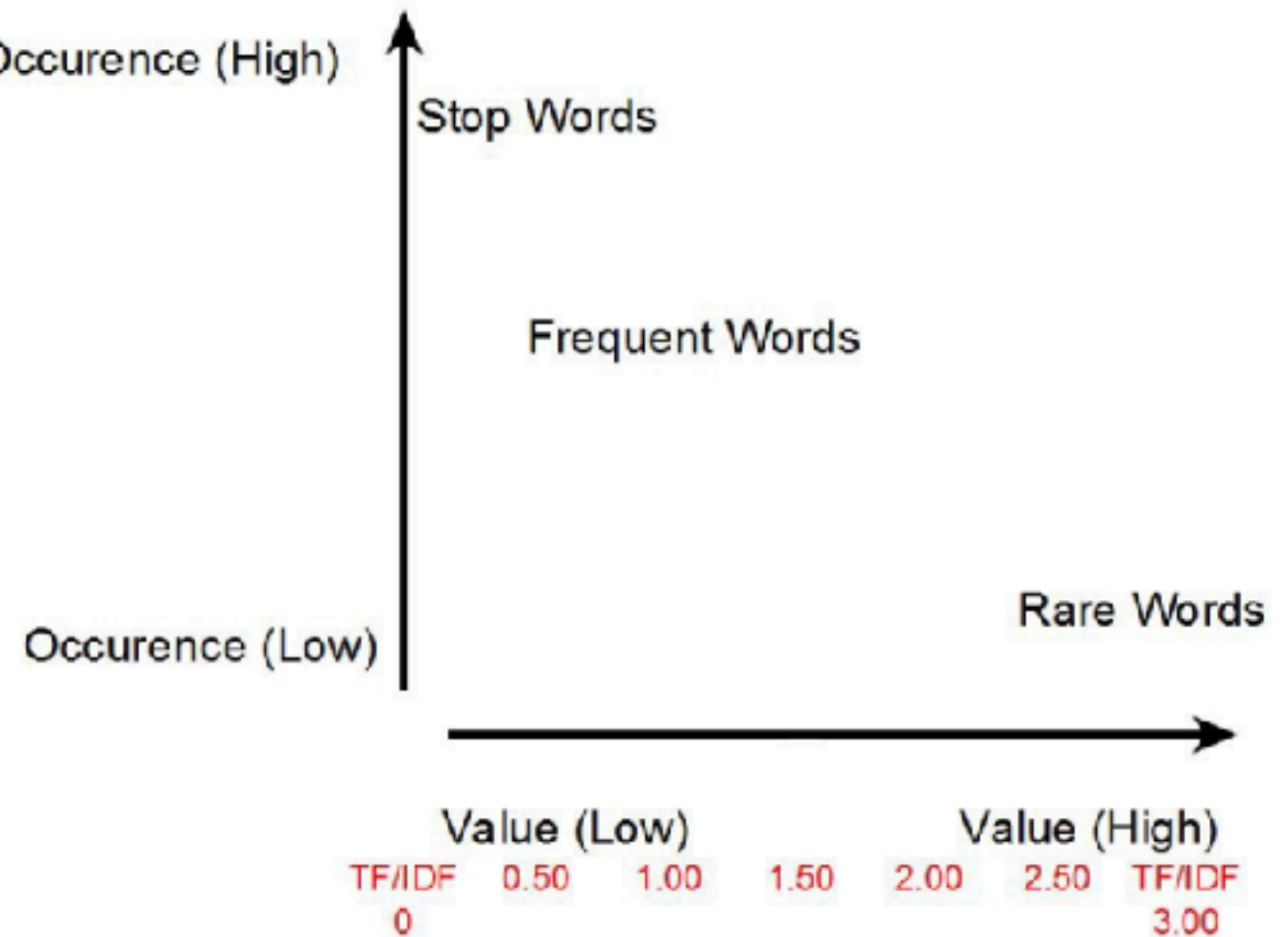
- Idea: Commonly occurring words don't carry useful information but rare words do

$$tfidf(term, document) = tf(term, document) \cdot idf(term)$$

$$tf(term, document) = \frac{n_i}{\sum_{k=1}^V n_k}$$

Occurrence (High)

$$idf(term) = \log \frac{N}{n_t}$$



TF-IDF Encoding

```
from sklearn.feature_extraction.text import TfidfTransformer

# create tf-idf object
transformer = TfidfTransformer(smooth_idf=False)

# X can be obtained as X.toarray() from the previous snippet
X = [[3, 0, 1],
      [5, 0, 0],
      [3, 0, 0],
      [1, 0, 0],
      [3, 2, 0],
      [3, 0, 4]]

# learn the vocabulary and store tf-idf sparse matrix in tfidf
tfidf = transformer.fit_transform(counts)

# retrieving matrix in numpy form as we did it before
tfidf.toarray()
```

Uninterested

prefix

stem

suffix

PREFIX/SUFFIX/ROOT

1. acou-
2. aer-, aero-
3. alt-, alti-, alto-
4. ampli-
5. angl, angul-
6. ann-, annu-
7. ap-, apo-
8. astr-, astro-, aster-
9. -ation, -ition, -tion
10. atmo-
11. -atude, -itude, -tude
12. aud-, audio-
13. bar-, baro-
14. calor-
15. can-, cand-
16. cap-, capac-
17. carb-, carbo-
18. -ced, -ceed, -cess
19. -celer
20. centr-, centri-
21. chem-, chemo-
22. chrom-, chromo-
23. chron-, chrono-
24. circ-, circl-
25. con-, com-
26. con-, coni-
27. cosm-, cosmo-
28. cry-, cryo-
29. cycly-, cyclo-
30. deca-
31. di-
32. -duc, -duce, -duct
33. dyn-, dynam-
34. elect-, electro-
35. -ence
36. end-, endo-, ento-
37. equa-, equi-, equ-
38. erg-, ergo-
39. ex-, exo-
40. ferro-

MEANING

- hear
- air
- height
- large
- angle
- year
- away from
- star
- process, result
- air
- condition, state of
- hear
- heavy
- heat
- glow white
- hold
- carbon
- going
- swift, to hasten
- center
- transmutation
- color
- time
- circle
- with, together
- cone
- universe
- cold
- circle
- ten
- two
- to lead
- power
- electric, amber
- state of
- inside, within
- equal
- work
- outside, out of
- iron

EXAMPLE

- acoustics
- aerodynamics
- altitude
- amplitude
- triangulation
- annual
- apogee
- astronomical
- condensation
- atmosphere
- amplitude
- audiologist
- barometric
- calorimetry
- incandescence
- capacitor
- carbohydrate
- recessional
- accelerometer
- centripetal
- chemical
- chromatograph
- chronograph
- circulation
- compound
- conical
- cosmological
- cryogenic
- cyclone
- dodecagon
- dipole
- aqueduct
- hydrodynamic
- electrician
- luminescence
- endothermic
- equilibrium
- ergonomics
- exogenic
- ferromagnetic

PREFIX/SUFFIX/ROOT

41. fiss-
42. -flect
43. flu-
44. for-, fort-
45. -fract
46. frict-
47. fus-
48. -fy
49. ge-, geo-, -gee
50. -gen
51. -gon
52. grad-
53. -gram
54. -graph
55. grav-
56. gyr-
57. helio-
58. hemi-
59. hepta-
60. hetero-
61. hexa-
62. homo-
63. hydr-, hydro-
64. -ic, -tic
65. -ical
66. -ician, -icist
67. -ics, -tics-
68. ign-
69. -ile
70. infra-
71. inter-
72. intra-
73. iso-
74. -ist
75. -istry
76. -ject
77. kin-, kine-
78. lept-
79. lev-
80. libr-, libri-

MEANING

- split
- bend
- flow
- strong, strength
- break
- rub
- melting
- to make
- earth
- to produce, origin
- angle
- step
- written record
- recording
- heavy, weighty
- rotate
- sun
- half
- seven
- unlike, different
- six
- same, like
- water
- person
- pertaining to
- specialist
- skill
- fire
- pertaining to, thing
- below
- between
- within
- equal
- person
- skill
- to throw
- motion, movement
- small
- to raise
- weight

EXAMPLE

- fission
- reflection
- flux
- fortification
- fracture
- friction
- fusion
- magnify
- geothermal
- hydrogen
- hexagonal
- gradient
- cardiogram
- spectrograph
- gravitational
- gyroscope
- heliocentric
- hemisphere
- heptathlon
- heterotroph
- hexagon
- homophone
- hydrothermal
- lunatic
- mechanical
- physicist
- mechanics
- ignition
- projectile
- infrasonic
- interstitial
- intramural
- isobaric
- chemist
- artistry
- projectile
- kinematics
- lepton
- levitation
- equilibrium

I love NLP and I like dogs

I = [_ _ _ _ _]

Love = [_ _ _ _ _]

NLP = [_ _ _ _ _]

And = [_ _ _ _ _]

Like = [_ _ _ _ _]

Dogs = [_ _ _ _ _]

Vector Filling Approaches

1. Word Count

2. Co-occurrence

and so on

I love NLP and I like dogs

	I	Love	NLP	And	Like	Dogs
I	0	1	0	1	1	0
Love	1	0	1	0	0	0
NLP	0	1	0	1	0	0
And	1	0	1	0	0	0
Like	1	0	0	0	0	1
Dogs	0	0	0	0	1	0

$$\mathbf{I} = [0 \ 1 \ 0 \ 1 \ 1 \ 0]$$

$$\mathbf{Love} = [1 \ 0 \ 1 \ 0 \ 0 \ 0]$$

$$\mathbf{NLP} = [0 \ 1 \ 0 \ 1 \ 0 \ 0]$$

$$\mathbf{And} = [1 \ 0 \ 1 \ 0 \ 0 \ 0]$$

$$\mathbf{Like} = [1 \ 0 \ 0 \ 0 \ 0 \ 1]$$

$$\mathbf{Dogs} = [0 \ 0 \ 0 \ 0 \ 1 \ 0]$$

Distributional Hypothesis

*Words that appear in same context
share semantic meaning*

1. Count-based method (Latent Semantic Analysis)

Count-based methods compute the statistics of how often some word co-occurs with its neighbour words in a large text corpus, and then map these count-statistics down to a small, dense vector for each word.

2. Predictive method (Neural Probabilistic Language Model)

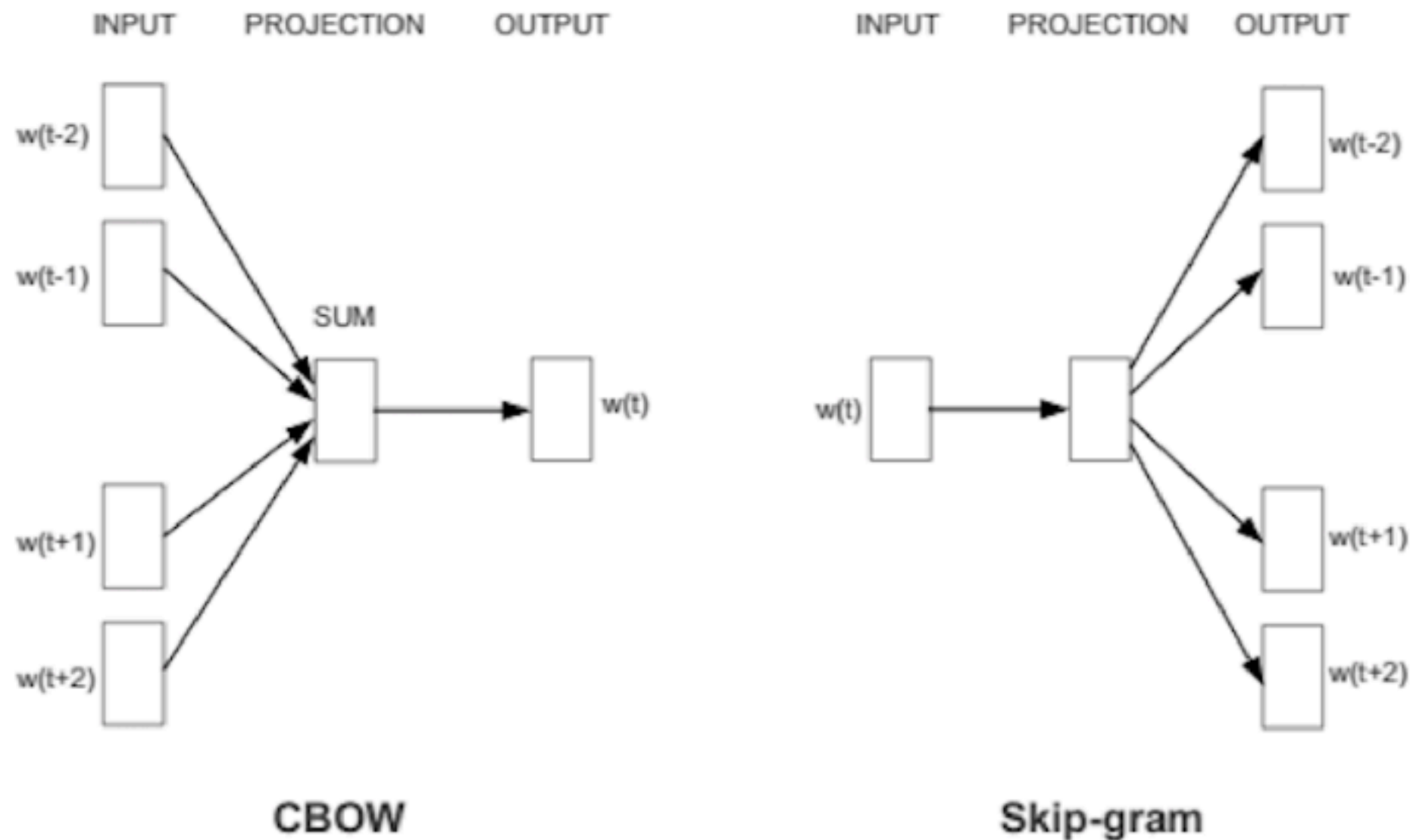
Predictive models directly try to predict a word from its neighbours in terms of learned small, dense embedding vectors.

Word2Vec —> Predictive Method

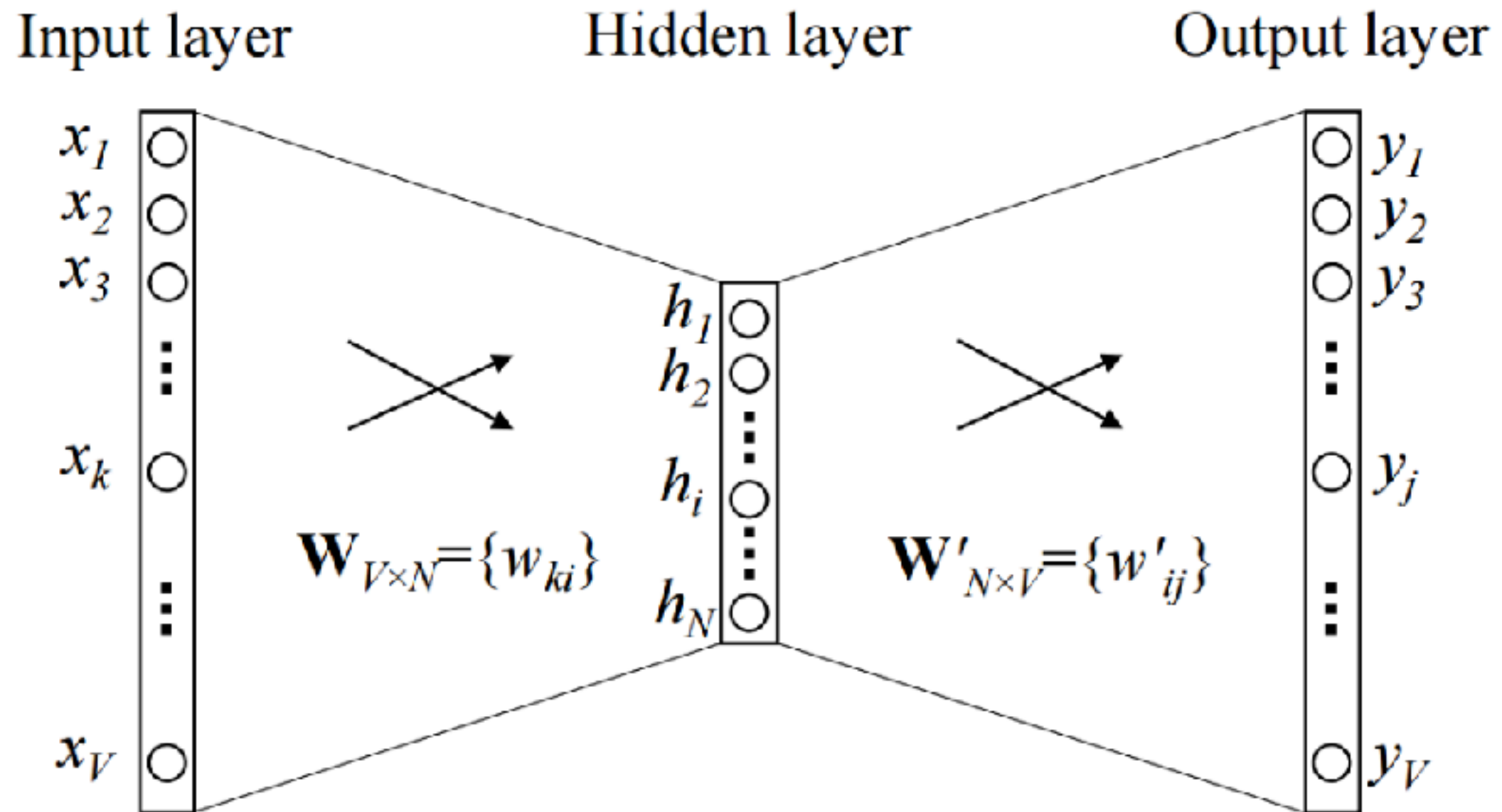
Word2Vec

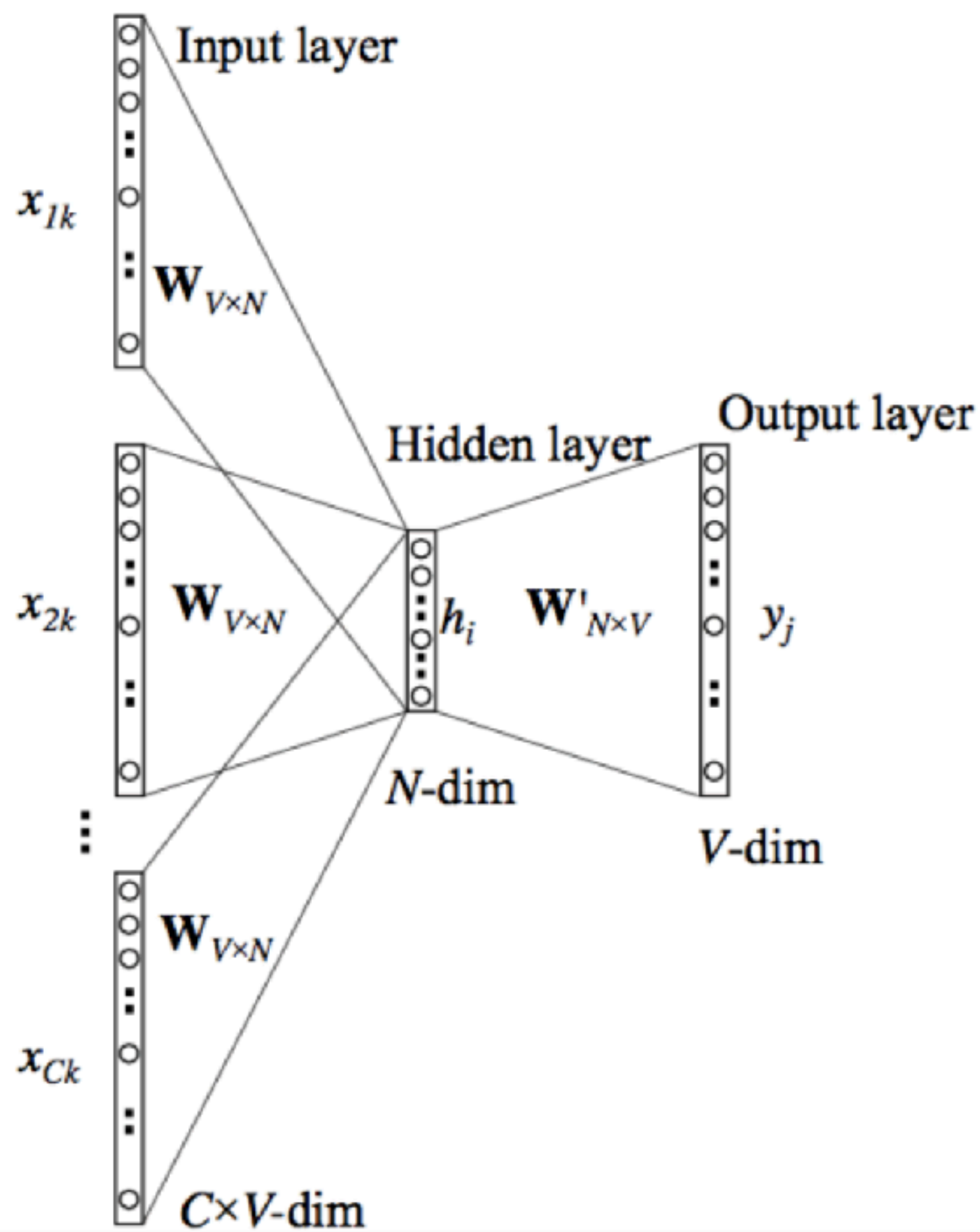
Word2vec is to group the vectors of similar words together in vectorspace. Word2vec creates vectors that are distributed numerical representations of word features, features such as the context of individual words.

Word2vec is similar to an autoencoder, encoding each word in a vector, but rather than training against the input words through reconstruction, as a [restricted Boltzmann machine](#) does, word2vec trains words against other words that neighbor them in the input corpus.

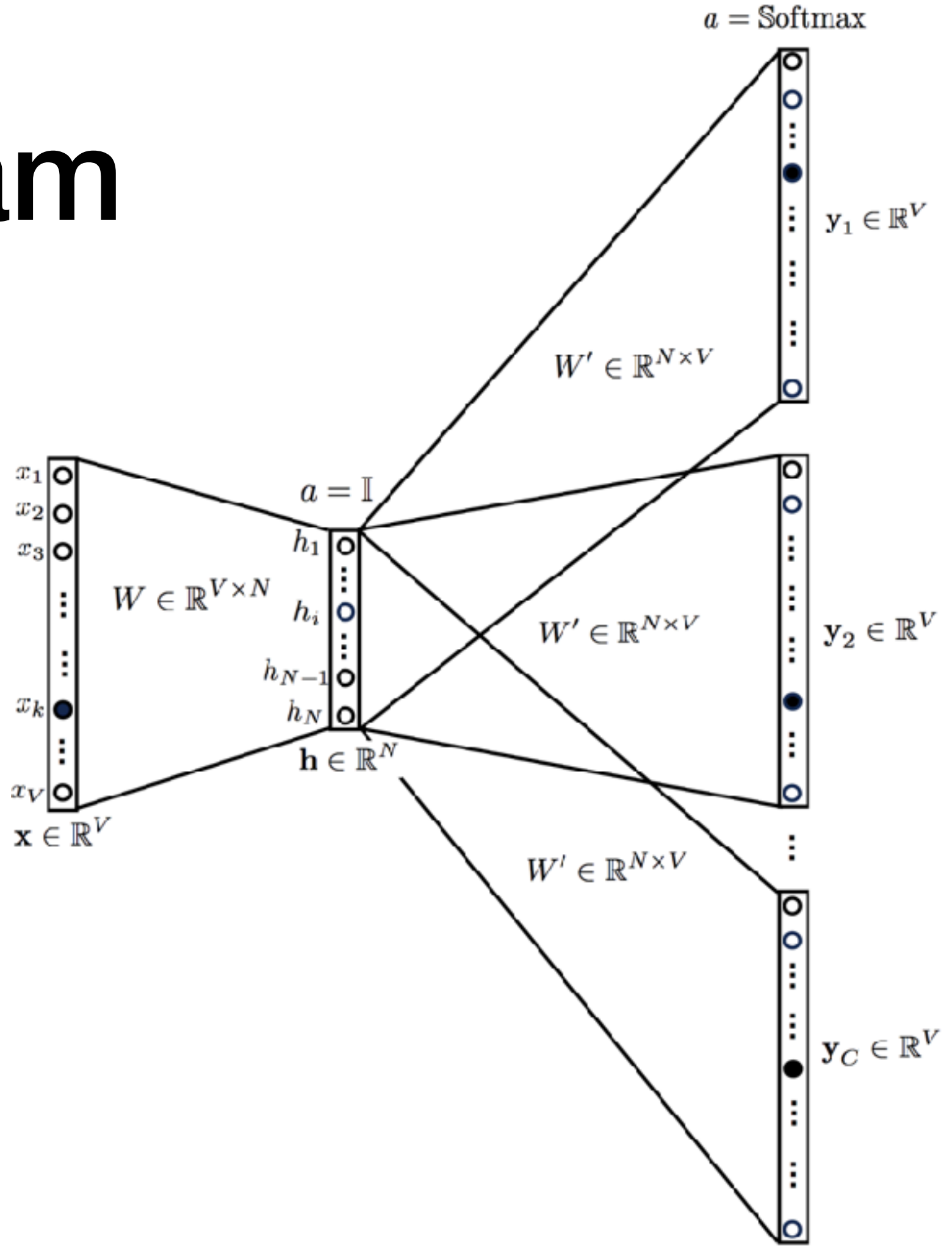


CBoW





Skip-Gram



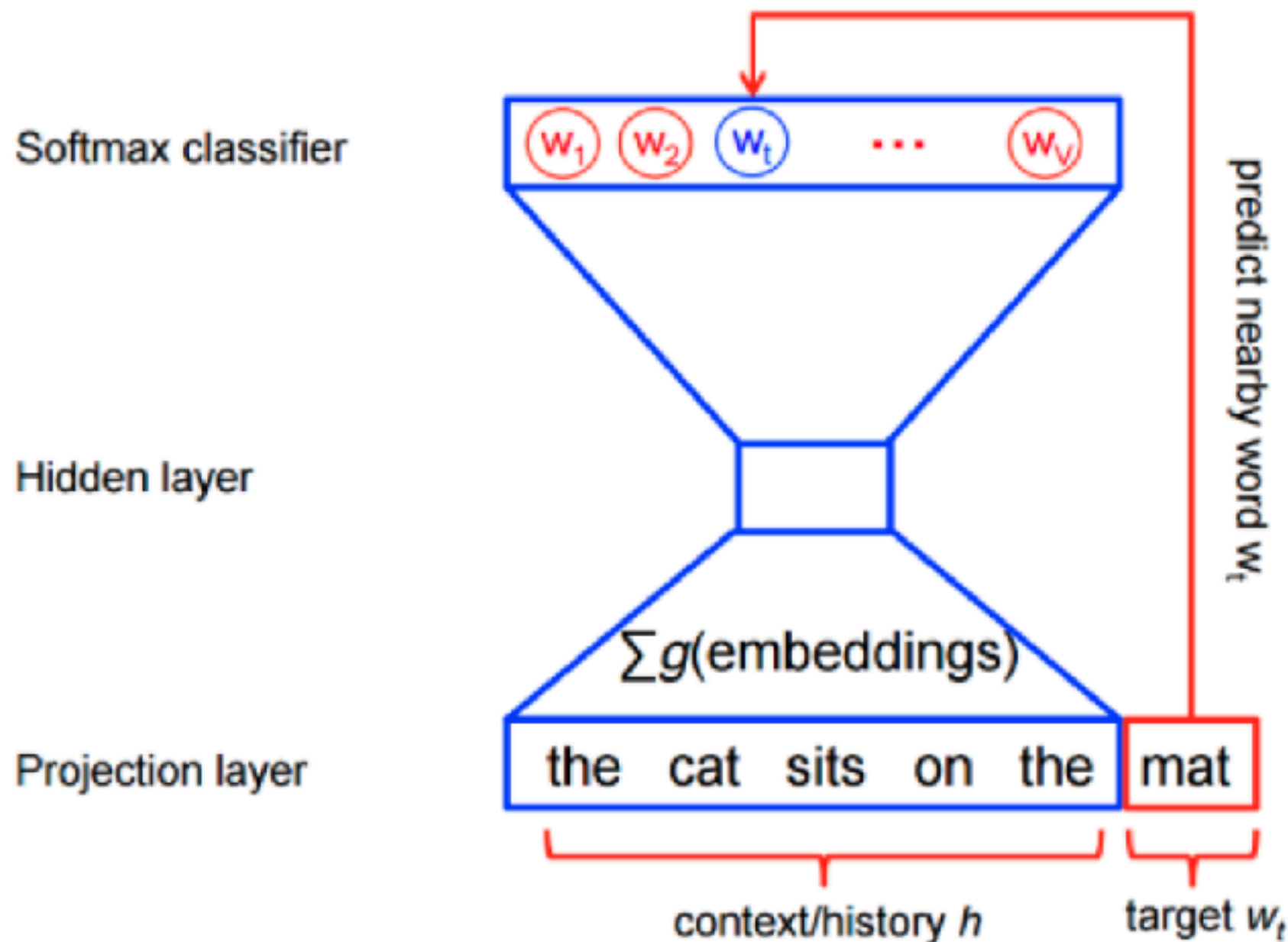
$$\begin{aligned}
 P(w_t|h) &= \text{softmax}(\text{score}(w_t, h)) \\
 &= \frac{\exp\{\text{score}(w_t, h)\}}{\sum_{\text{Word } w' \text{ in Vocab}} \exp\{\text{score}(w', h)\}}
 \end{aligned}$$

$$\begin{aligned}
 J_{\text{ML}} &= \log P(w_t|h) \\
 &= \text{score}(w_t, h) - \log \left(\sum_{\text{Word } w' \text{ in Vocab}} \exp\{\text{score}(w', h)\} \right)
 \end{aligned}$$

$$\text{argmax}_{\theta} \frac{1}{T} \sum_{t=1}^T \sum_{j \in c, j \neq 0} \log p(w_{t+j}|w_t; \theta)$$

Embedding Lookup Matrix

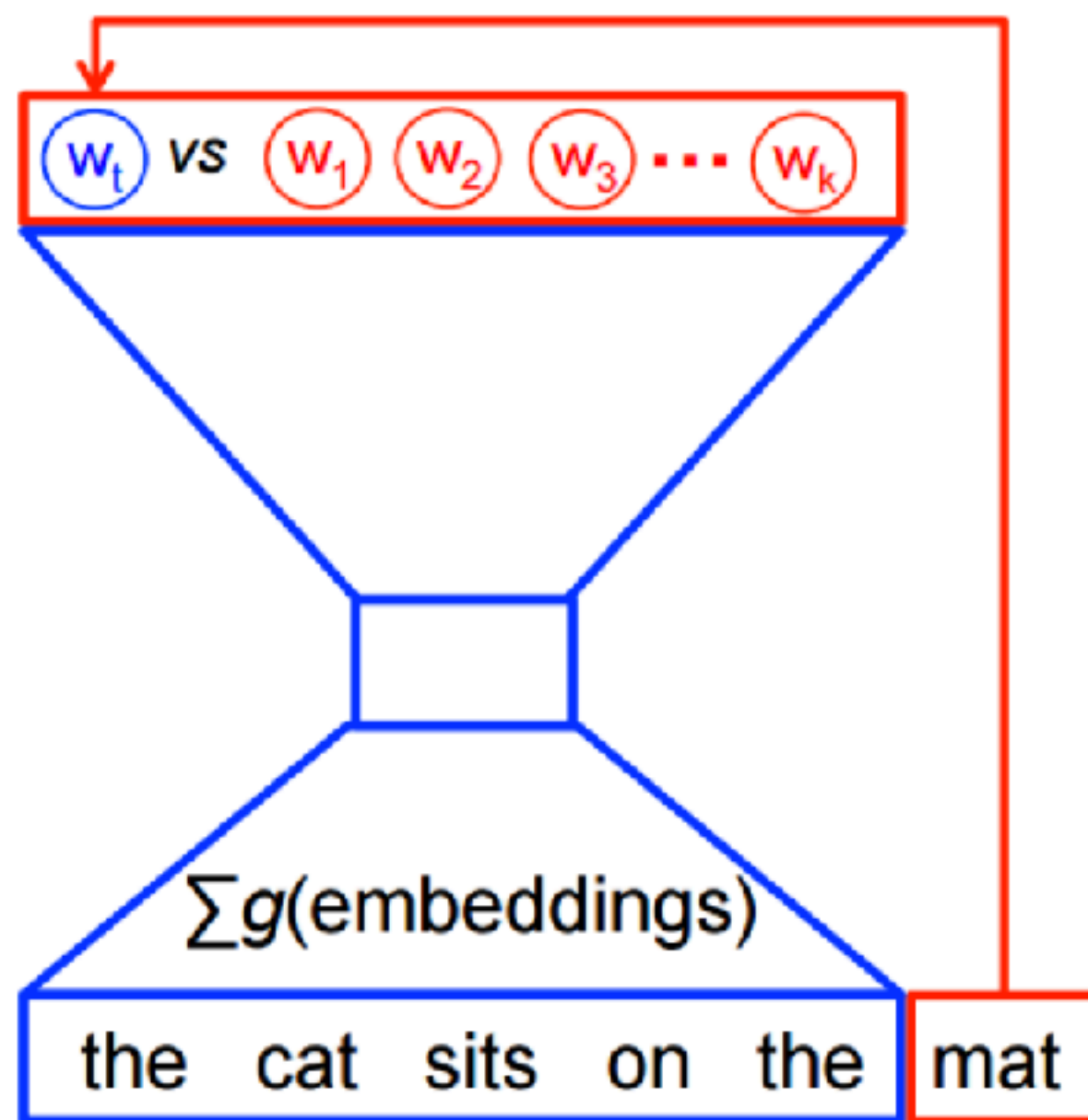
$$p(w_i|w_t; \theta) = \frac{\exp(\theta w_i)}{\sum_t \exp(\theta w_t)}$$



Noise classifier

Hidden layer

Projection layer



Negative Sampling

$$J_{\text{NEG}} = \log Q_{\theta}(D = 1|w_t, h) + k \mathbb{E}_{\tilde{w} \sim P_{\text{noise}}} [\log Q_{\theta}(D = 0|\tilde{w}, h)]$$

where $Q_{\theta}(D = 1|w, h)$ is the binary logistic regression probability under the model of seeing the word w in the context h in the dataset D , calculated in terms of the learned embedding vectors θ . In practice we approximate the expectation by drawing k contrastive words from the noise distribution (i.e. we compute a [Monte Carlo average](#)).

$$J_{\text{NEG}}^{(t)} = \log Q_{\theta}(D = 1|\text{the, quick}) + \log(Q_{\theta}(D = 0|\text{sheep, quick}))$$

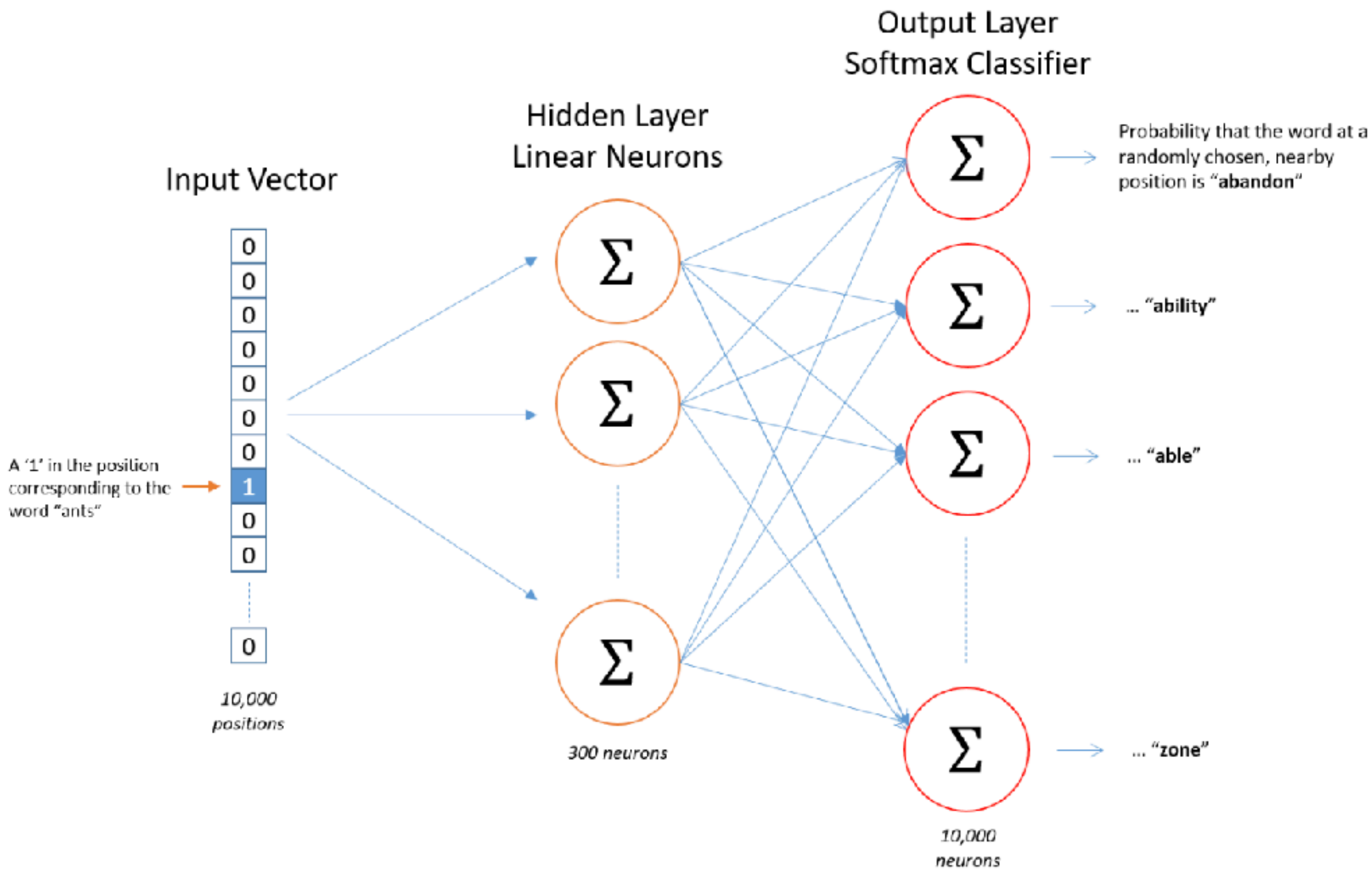
Word2Vec Steps

1. Take a 3 layer neural network. (1 input layer + 1 hidden layer + 1 output layer)
2. Feed it a word and train it to predict its neighbouring word.
3. Remove the last (output layer) and keep the input and hidden layer.
4. Now, input a word from within the vocabulary. The output given at the hidden layer is the '*word embedding*' of the input word.

Source Text

Training Samples

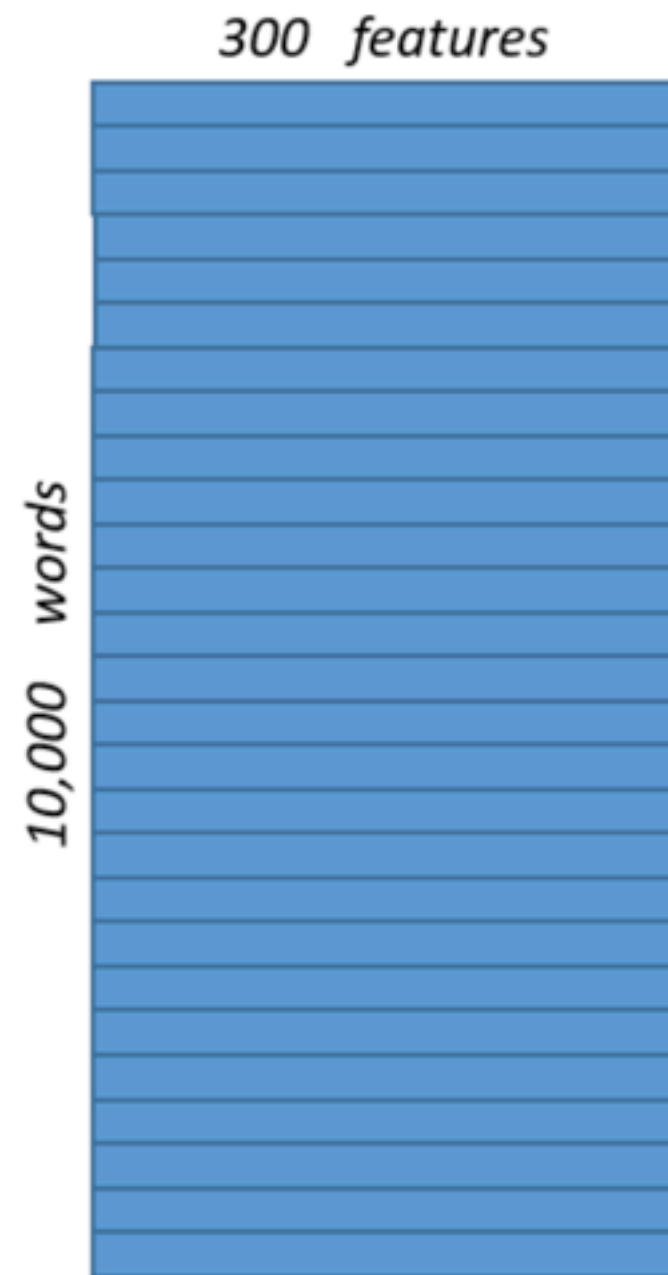
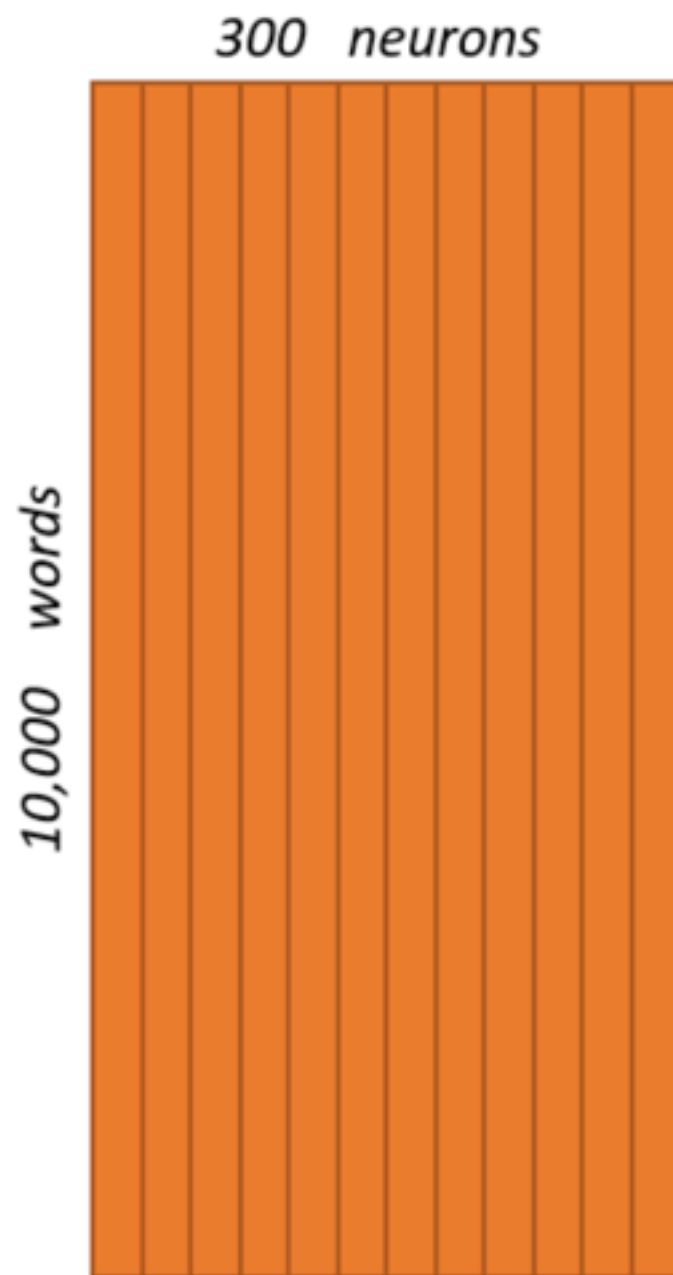
The quick brown fox jumps over the lazy dog. →	(the, quick) (the, brown)
The quick brown fox jumps over the lazy dog. →	(quick, the) (quick, brown) (quick, fox)
The quick brown fox jumps over the lazy dog. →	(brown, the) (brown, quick) (brown, fox) (brown, jumps)
The quick brown fox jumps over the lazy dog. →	(fox, quick) (fox, brown) (fox, jumps) (fox, over)



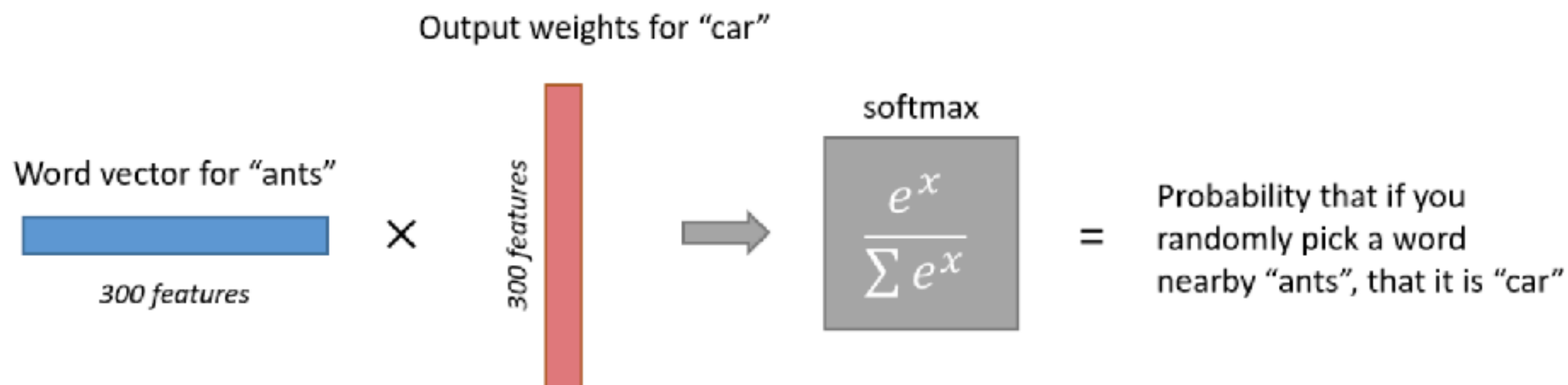
Hidden Layer
Weight Matrix



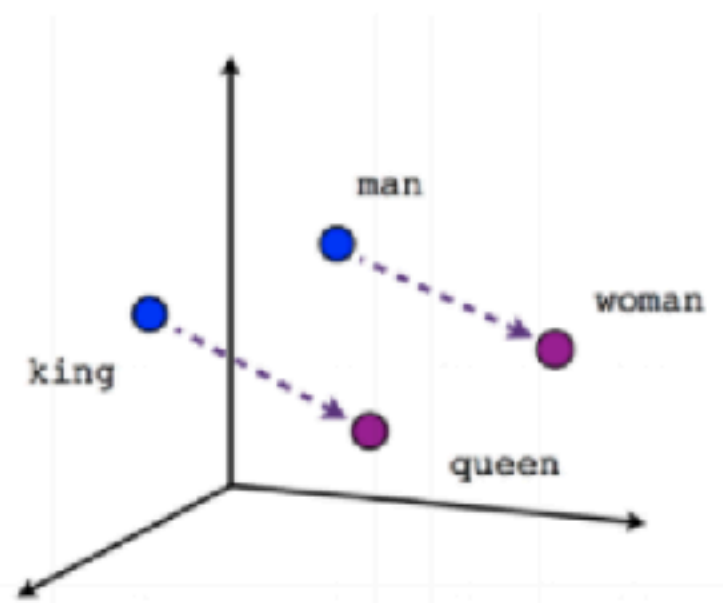
*Word Vector
Lookup Table!*



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



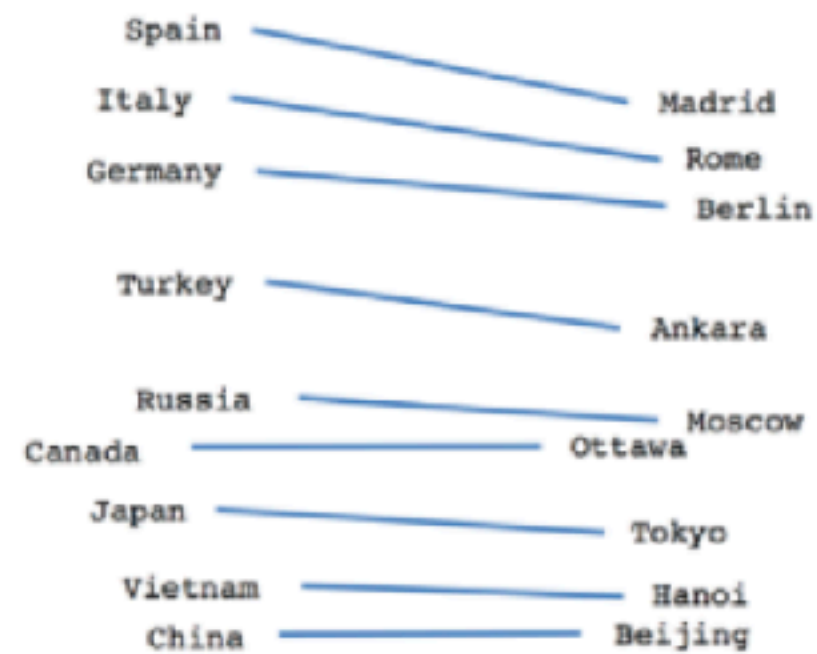
Word2Vec Relationships



Male-Female



Verb tense



Country-Capital

Implementation in Tensorflow

```
embeddings = tf.Variable(  
    tf.random_uniform([vocabulary_size, embedding_size], -1.0, 1.0))
```

```
nce_weights = tf.Variable(  
    tf.truncated_normal([vocabulary_size, embedding_size],  
                        stddev=1.0 / math.sqrt(embedding_size)))  
nce_biases = tf.Variable(tf.zeros([vocabulary_size]))
```

```
# Placeholders for inputs
```

```
train_inputs = tf.placeholder(tf.int32, shape=[batch_size])
```

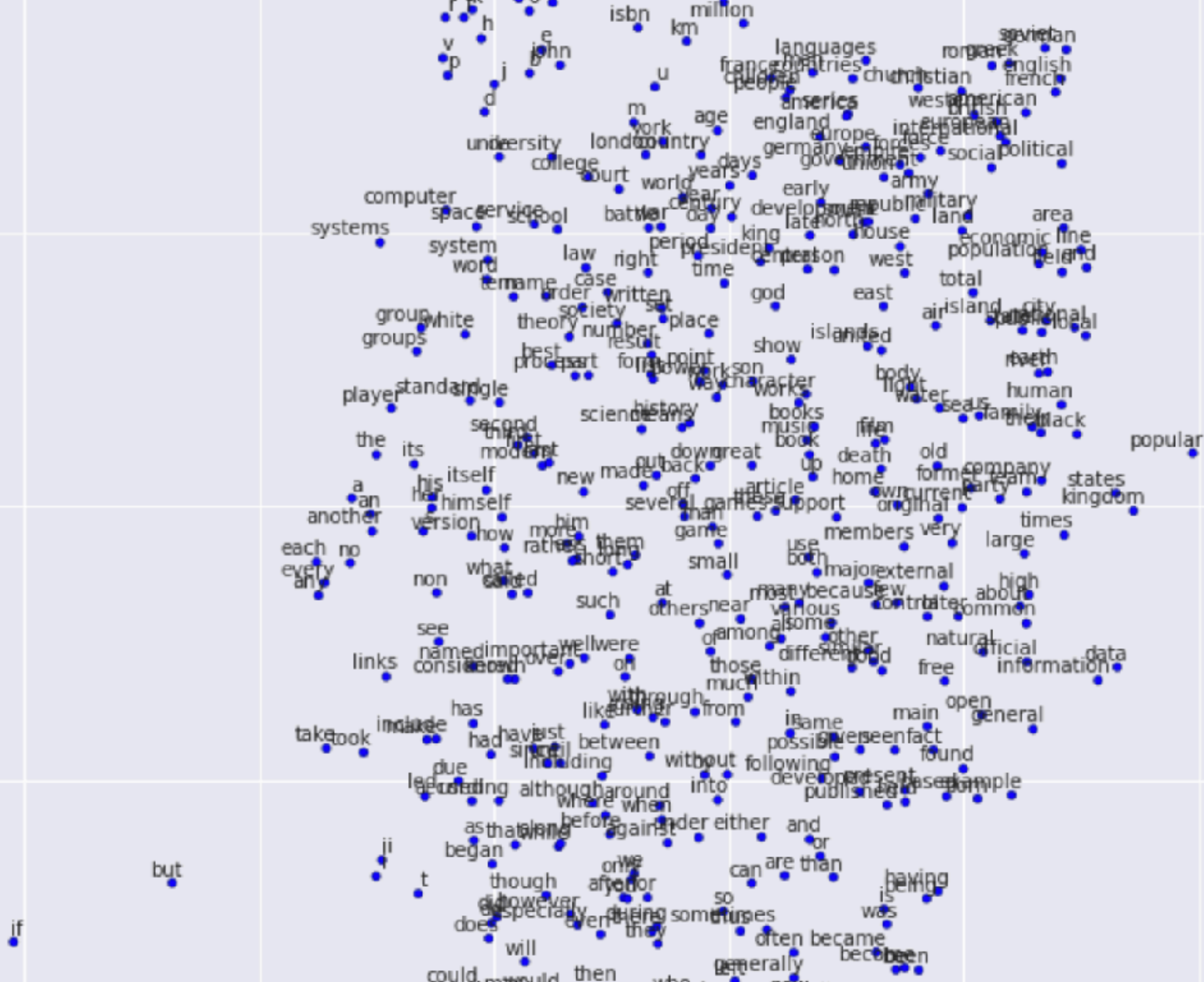
```
train_labels = tf.placeholder(tf.int32, shape=[batch_size, 1])
```

```
embed = tf.nn.embedding_lookup(embeddings, train_inputs)
```

```
# Compute the NCE loss, using a sample of the negative labels each time.
loss = tf.reduce_mean(
    tf.nn.nce_loss(weights=nce_weights,
                   biases=nce_biases,
                   labels=train_labels,
                   inputs=embed,
                   num_sampled=num_sampled,
                   num_classes=vocabulary_size))
```

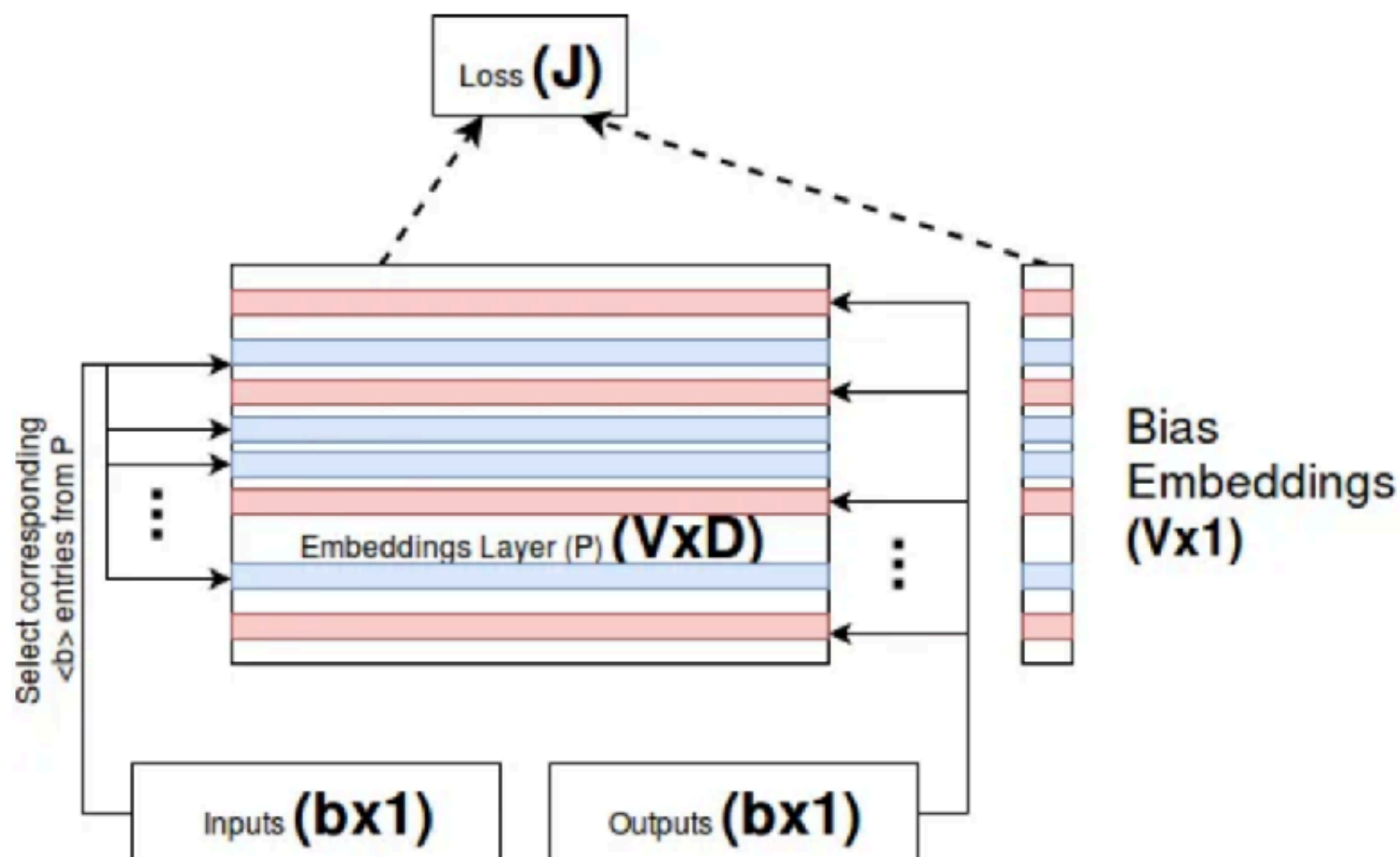
```
# We use the SGD optimizer.
optimizer = tf.train.GradientDescentOptimizer(learning_rate=1.0).minimize(loss)
```

```
for inputs, labels in generate_batch(...):  
    feed_dict = {train_inputs: inputs, train_labels: labels}  
    _, cur_loss = session.run([optimizer, loss], feed_dict=feed_dict)
```

GloVe

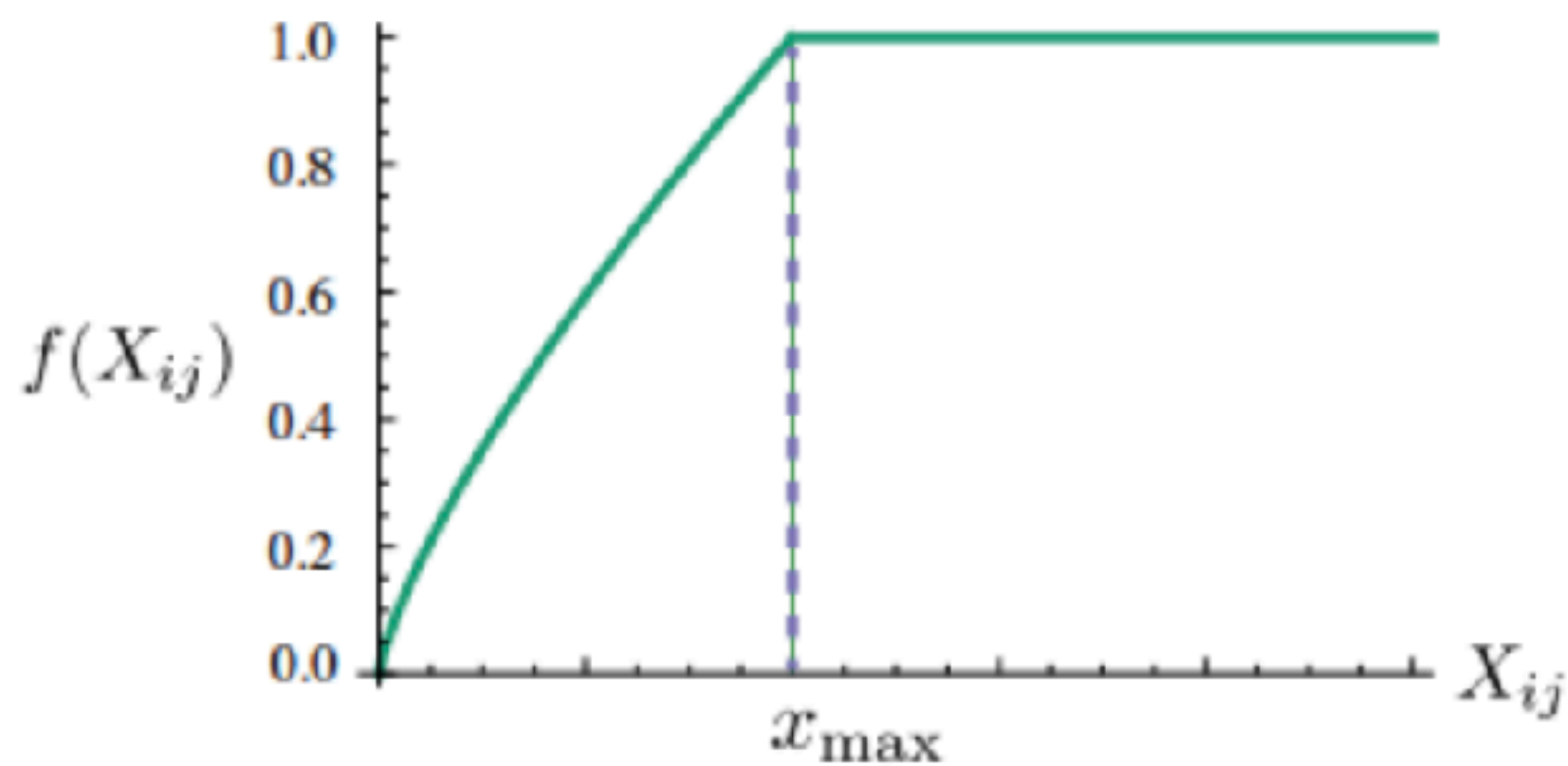
- Idea: Ratios between probabilities of words appearing next to each other carry more information than individual probabilities
 1. $P_{ice,solid} / P_{steam,solid}$ will be very high
 2. $P_{ice,gas} / P_{steam,gas}$ is very low
 3. $P_{ice,water} / P_{steam,water}$ will be higher than $P_{ice,fashion} / P_{steam,fashion}$



GloVe: High-level Architecture

Probability and Ratio	$k = solid$	$k = gas$	$k = water$	$k = fashion$
$P(k ice)$	1.9×10^{-4}	6.6×10^{-5}	3.0×10^{-3}	1.7×10^{-5}
$P(k steam)$	2.2×10^{-5}	7.8×10^{-4}	2.2×10^{-3}	1.8×10^{-5}
$P(k ice)/P(k steam)$	8.9	8.5×10^{-2}	1.36	0.96

$$f(x) = \begin{cases} (x/x_{\max})^\alpha & \text{if } x < x_{\max} \\ 1 & \text{otherwise} \end{cases}$$



Word Analogy

Model	Dim.	Size	Sem.	Syn.	Tot.
ivLBL	100	1.5B	55.9	50.1	53.2
HPCA	100	1.6B	4.2	16.4	10.8
GloVe	100	1.6B	<u>67.5</u>	<u>54.3</u>	<u>60.3</u>
SG	300	1B	61	61	61
CBOW	300	1.6B	16.1	52.6	36.1
vLBL	300	1.5B	54.2	<u>64.8</u>	60.0
ivLBL	300	1.5B	65.2	63.0	64.0
GloVe	300	1.6B	<u>80.8</u>	61.5	<u>70.3</u>
SVD	300	6B	6.3	8.1	7.3
SVD-S	300	6B	36.7	46.6	42.1
SVD-L	300	6B	56.6	63.0	60.1
CBOW [†]	300	6B	63.6	<u>67.4</u>	65.7
SG [†]	300	6B	73.0	66.0	69.1
GloVe	300	6B	<u>77.4</u>	67.0	<u>71.7</u>
CBOW	1000	6B	57.3	68.9	63.7
SG	1000	6B	66.1	65.1	65.6
SVD-L	300	42B	38.4	58.2	49.2
GloVe	300	42B	<u>81.9</u>	<u>69.3</u>	<u>75.0</u>

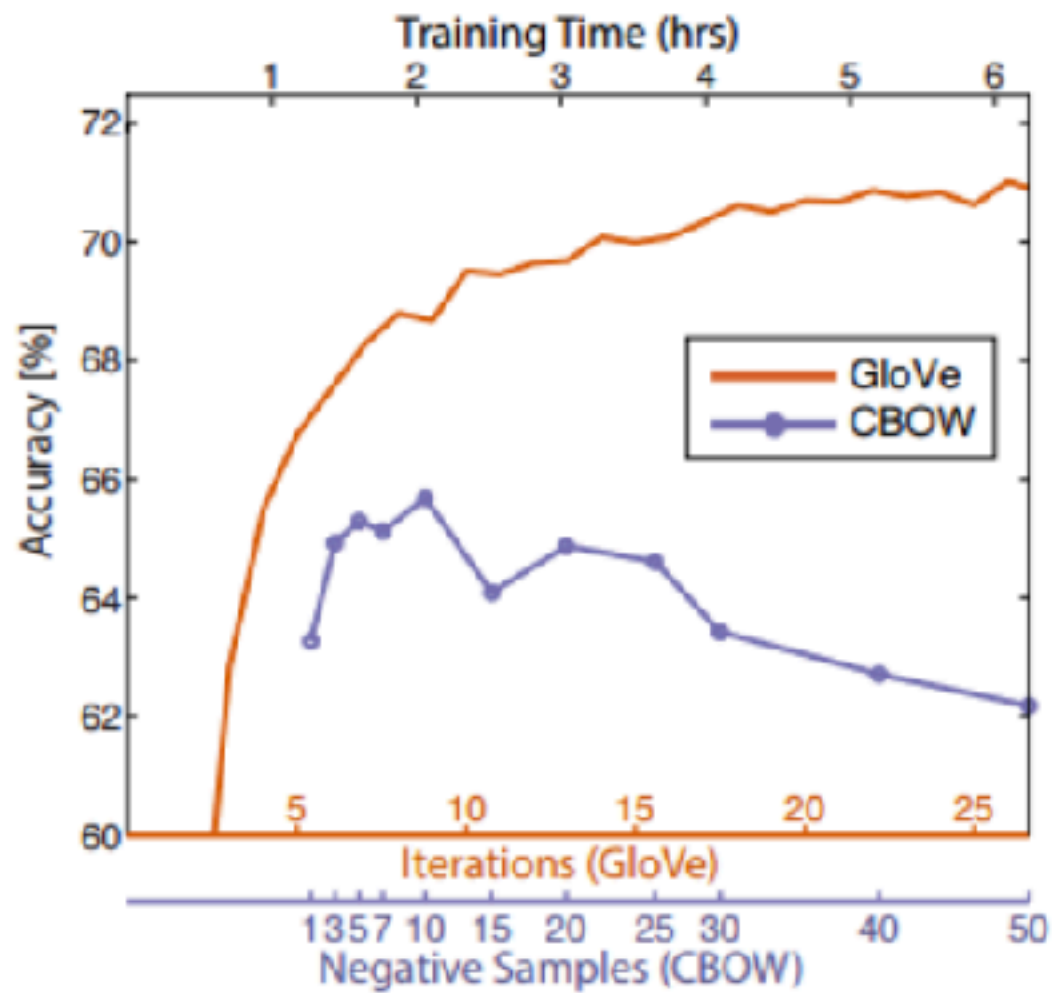
Word Similarity

Model	Size	WS353	MC	RG	SCWS	RW
SVD	6B	35.3	35.1	42.5	38.3	25.6
SVD-S	6B	56.5	71.5	71.0	53.6	34.7
SVD-L	6B	65.7	<u>72.7</u>	75.1	56.5	37.0
CBOW [†]	6B	57.2	65.6	68.2	57.0	32.5
SG [†]	6B	62.8	65.2	69.7	<u>58.1</u>	37.2
GloVe	6B	<u>65.8</u>	<u>72.7</u>	<u>77.8</u>	53.9	<u>38.1</u>
SVD-L	42B	74.0	76.4	74.1	58.3	39.9
GloVe	42B	<u>75.9</u>	<u>83.6</u>	<u>82.9</u>	<u>59.6</u>	<u>47.8</u>
CBOW [*]	100B	68.4	79.6	75.4	59.4	45.5

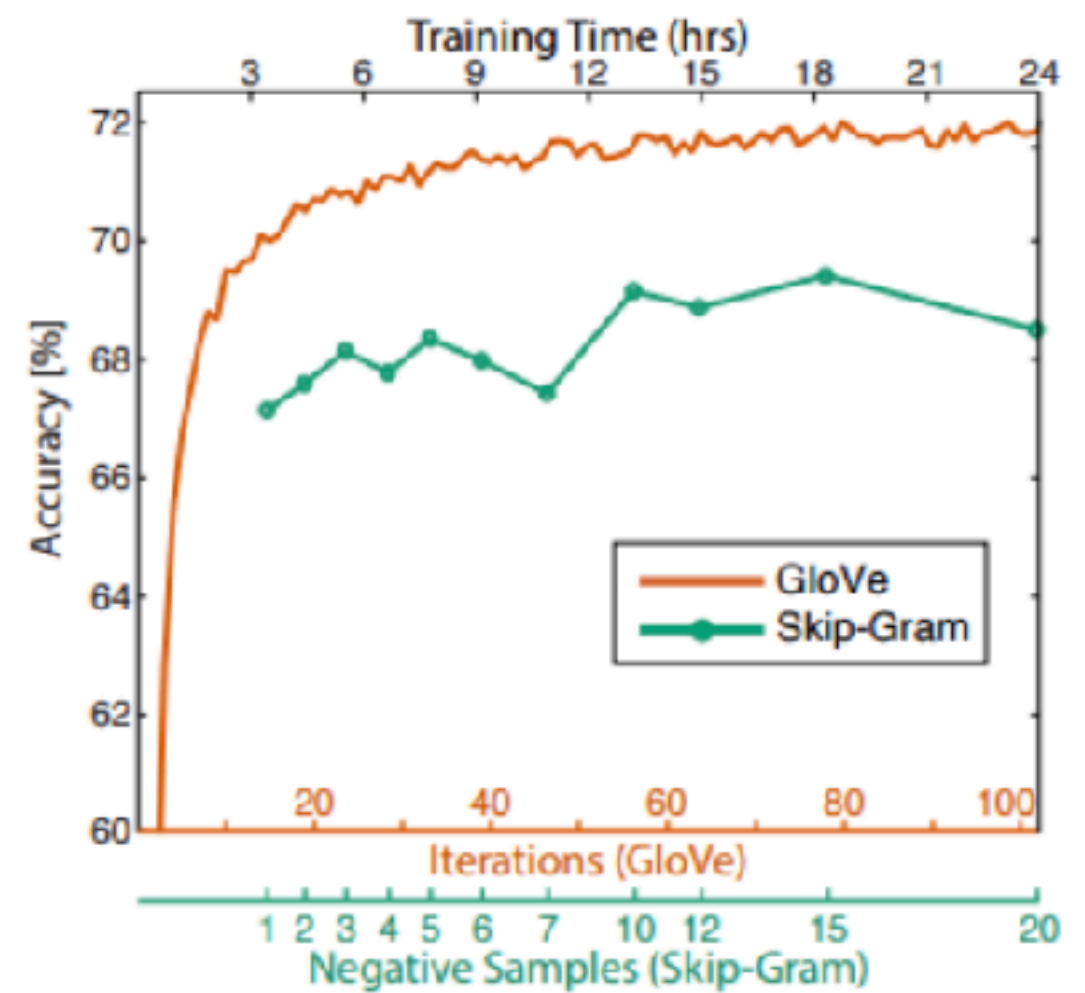
NER

Model	Dev	Test	ACE	MUC7
Discrete	91.0	85.4	77.4	73.4
SVD	90.8	85.7	77.3	73.7
SVD-S	91.0	85.5	77.6	74.3
SVD-L	90.5	84.8	73.6	71.5
HPCA	92.6	88.7	81.7	80.7
HSMN	90.5	85.7	78.7	74.7
CW	92.2	87.4	81.7	80.2
CBOW	93.1	88.2	82.2	81.1
GloVe	93.2	88.3	82.9	82.2

GloVe Vs Word2Vec



(a) GloVe vs CBOW



(b) GloVe vs Skip-Gram