



# Word Embedding

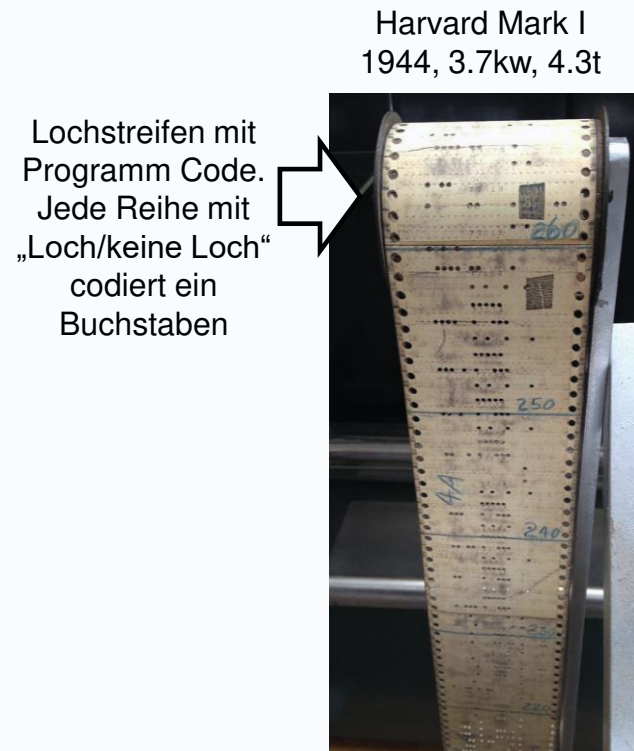
*From Punched Tape to Character encoding and Text-Preprocessing*

Used to represent a repertoire of characters by an encoding system that assigns a number to each character for digital representation, e.g. UTF-8, Unicode, ...

# Word Embedding

## From Punched Tape to Character encoding and Text-Preprocessing

Used to represent a repertoire of characters by an encoding system that assigns a number to each character for digital representation, e.g. UTF-8, Unicode, ...

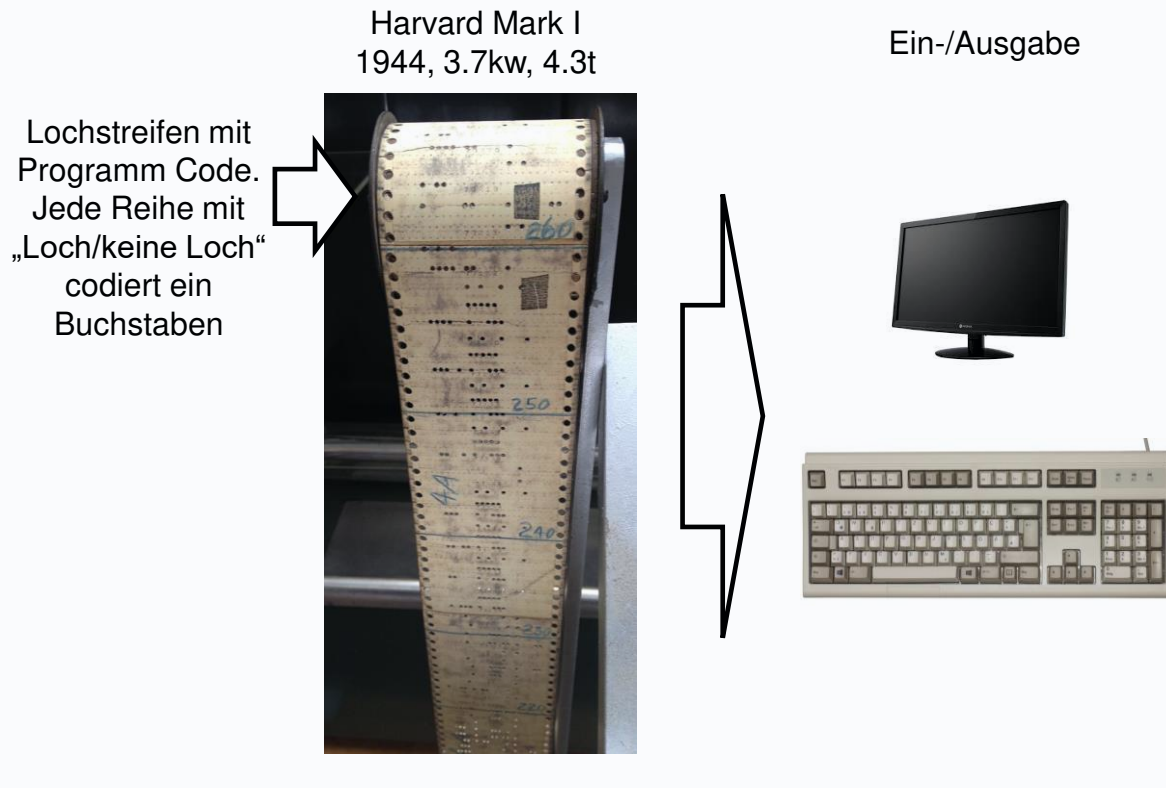


# Word Embedding

## From Punched Tape to Character encoding and Text-Preprocessing



Used to represent a repertoire of characters by an encoding system that assigns a number to each character for digital representation, e.g. UTF-8, Unicode, ...

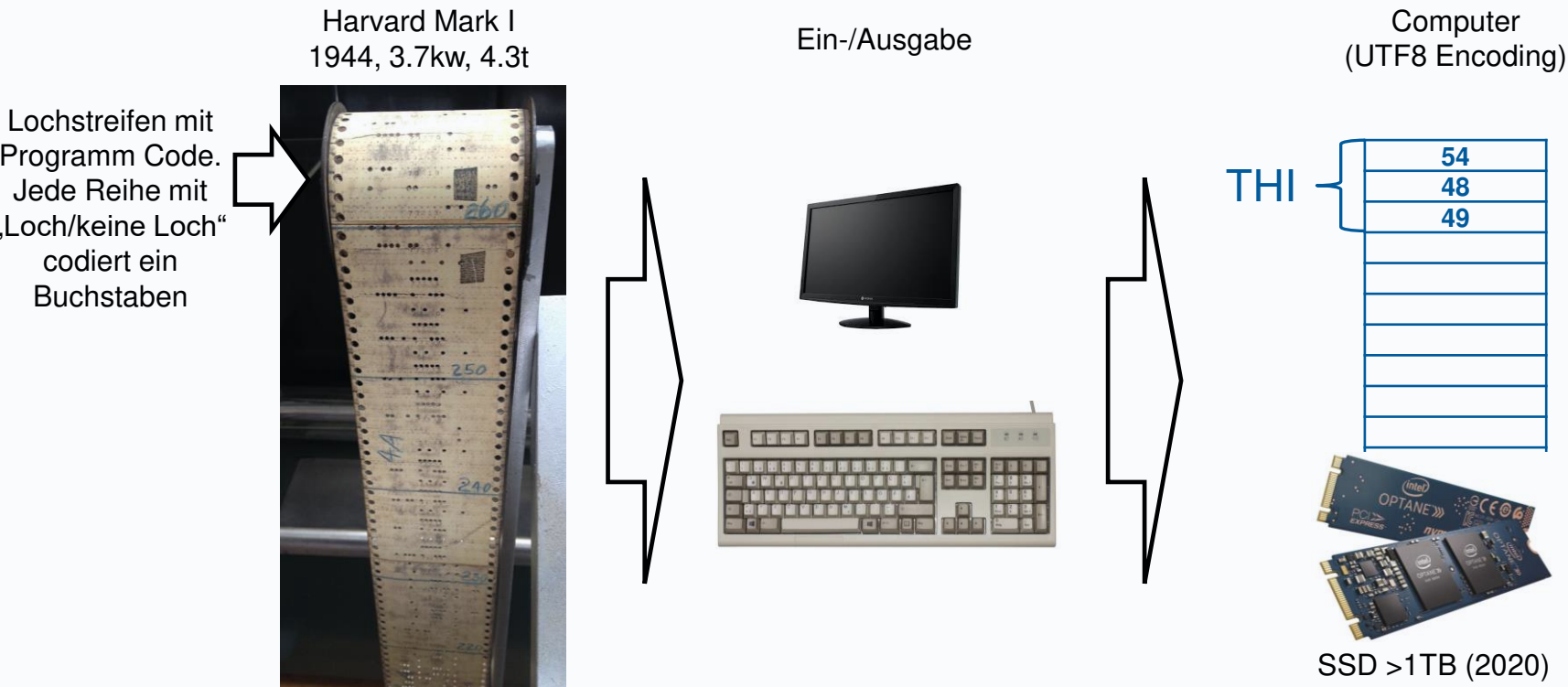


# Word Embedding

From Punched Tape to Character encoding and Text-Preprocessing



Used to represent a repertoire of characters by an encoding system that assigns a number to each character for digital representation, e.g. UTF-8, Unicode, ...

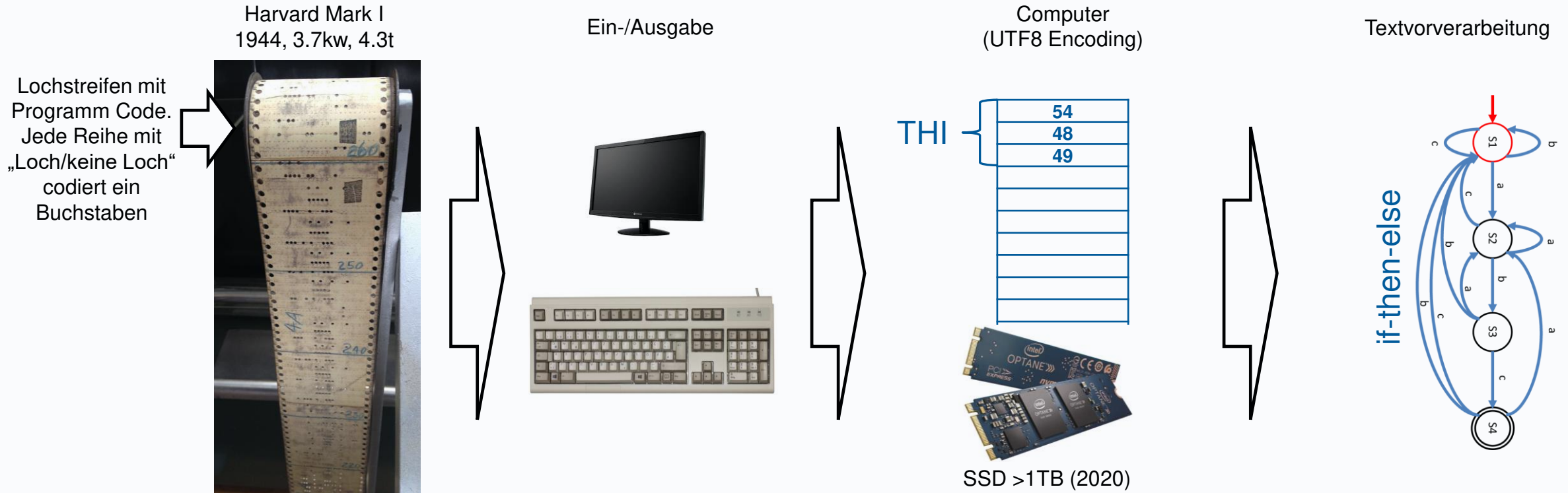


# Word Embedding

## From Punched Tape to Character encoding and Text-Preprocessing



Used to represent a repertoire of characters by an encoding system that assigns a number to each character for digital representation, e.g. UTF-8, Unicode, ...



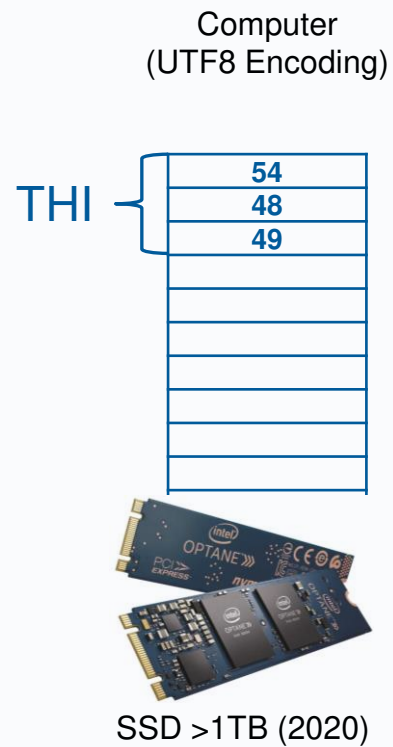


The data is encoded by a sequence of numbers. Let's try to interpret these numbers, mathematically.

# Word Embedding

## Mathematical view of a Character Encoding

The data is encoded by a sequence of numbers. Let's try to interpret these numbers, mathematically.

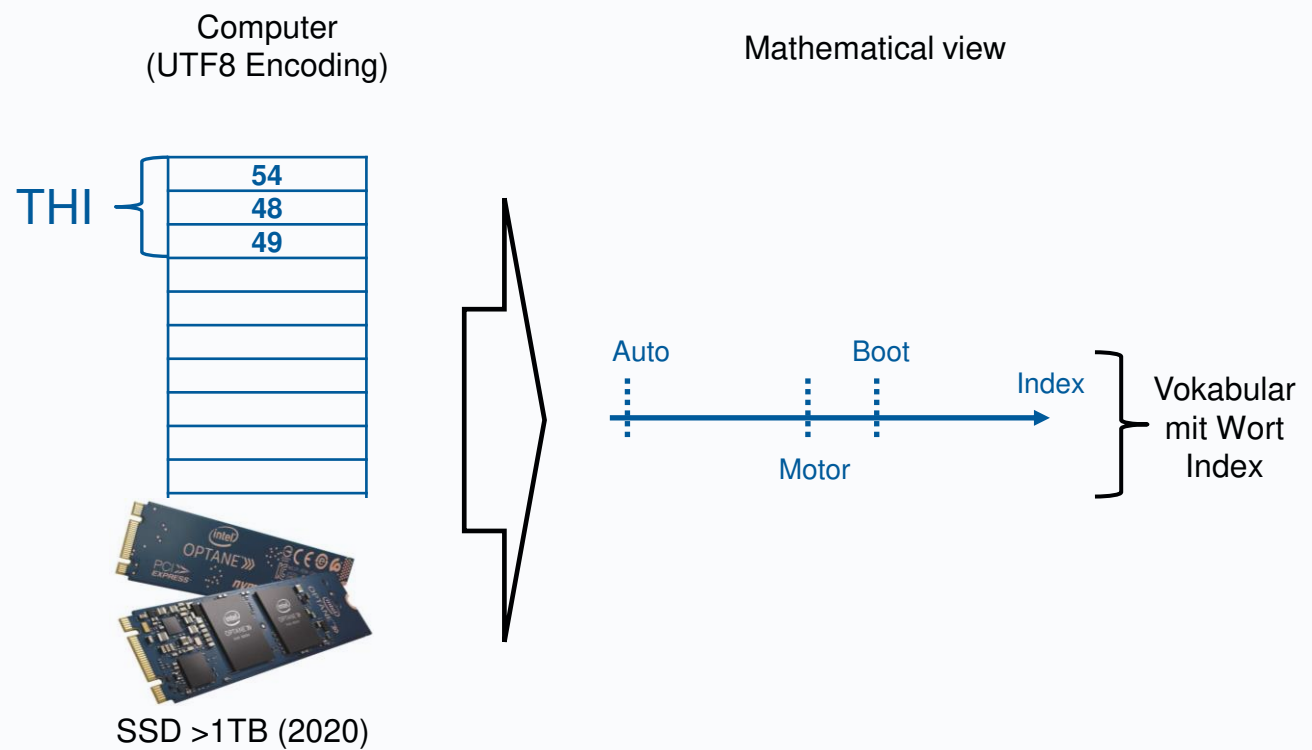


# Word Embedding

## Mathematical view of a Character Encoding



The data is encoded by a sequence of numbers. Let's try to interpret these numbers, mathematically.



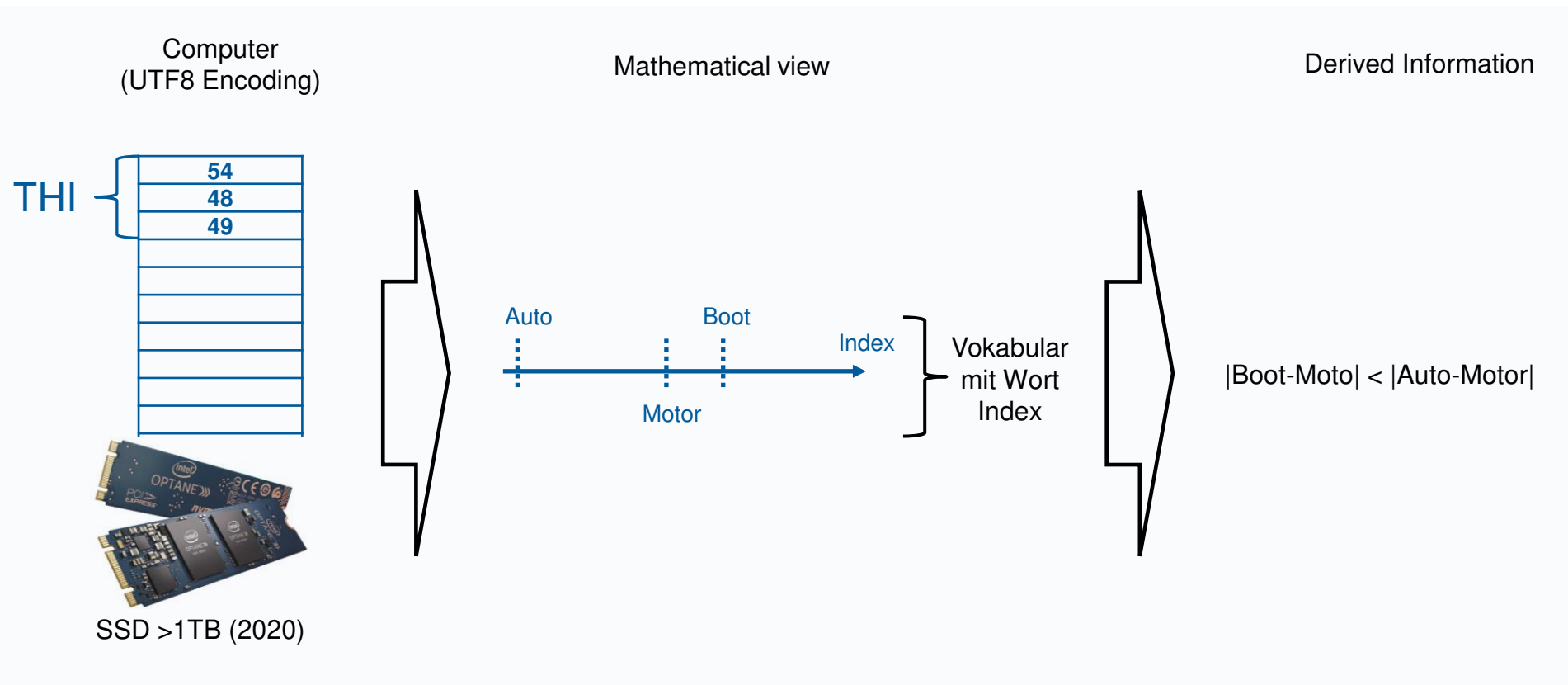


# Word Embedding

## Mathematical view of a Character Encoding



The data is encoded by a sequence of numbers. Let's try to interpret these numbers, mathematically.

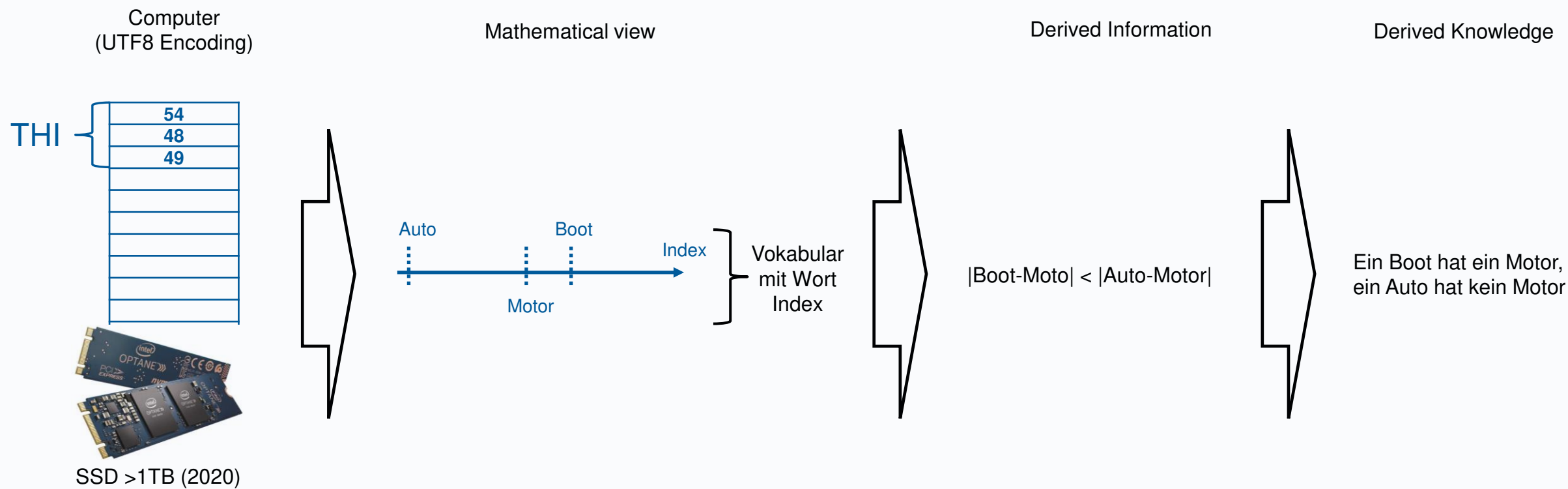


# Word Embedding

## Mathematical view of a Character Encoding



The data is encoded by a sequence of numbers. Let's try to interpret these numbers, mathematically.

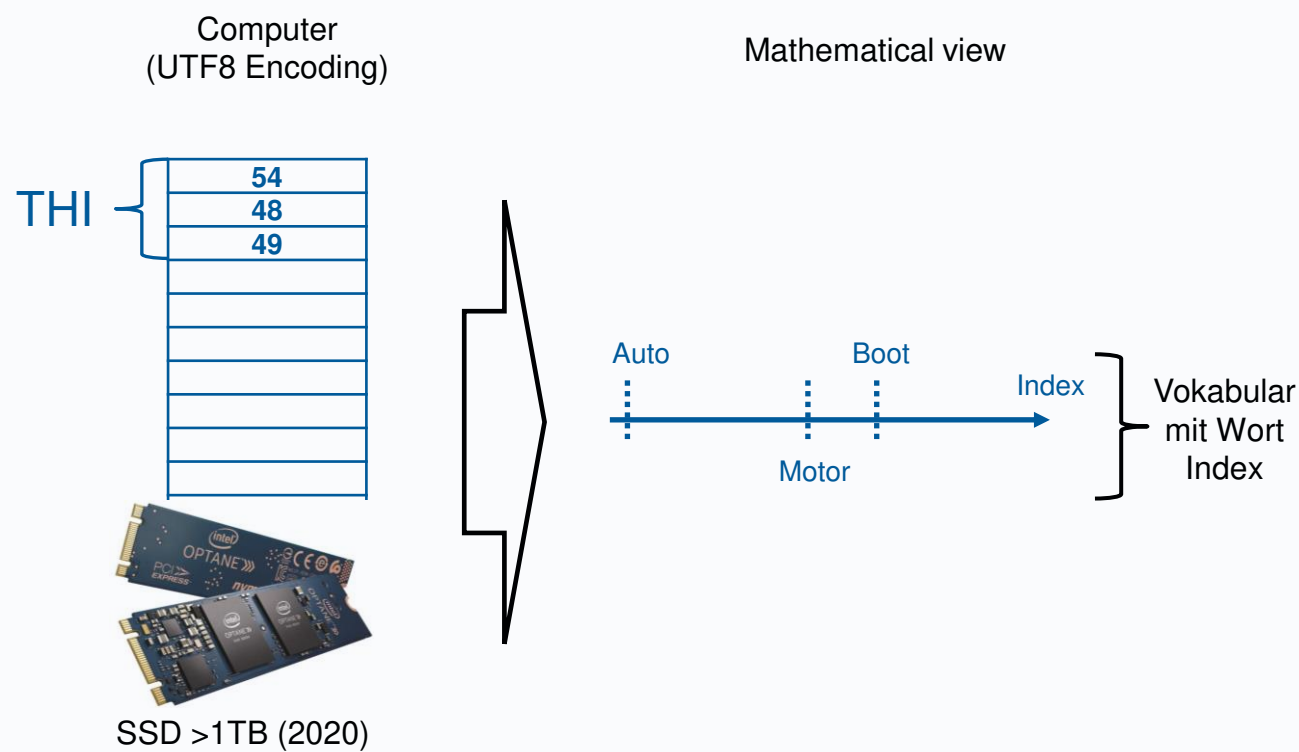


# Word Embedding

## Mathematical view of a Character Encoding



The data is encoded by a sequence of numbers. Let's try to interpret these numbers, mathematically.



Derived Information

Derived Knowledge

$|Boot-Moto| < |Auto-Motor|$

Ein Boot hat ein Motor, ein Auto hat kein Motor

**Doesn't work. No meaningful semantics is deducable**

# Word Embedding

*From storing characters to storing relational data*



Data using a structure and language consistent with first-order predicate logic, first described in 1969 by Edgar F. Codd, where all data is represented in terms of tuples, grouped into relations.

# Word Embedding

*From storing characters to storing relational data*

Data using a structure and language consistent with first-order predicate logic, first described in 1969 by Edgar F. Codd, where all data is represented in terms of tuples, grouped into relations.

Computer  
(UTF8 Encoding)

54
48
49



SSD >1TB (2020)

# Word Embedding

From storing characters to storing relational data



Data using a structure and language consistent with first-order predicate logic, first described in 1969 by Edgar F. Codd, where all data is represented in terms of tuples, grouped into relations.

Computer  
(UTF8 Encoding)

54
48
49



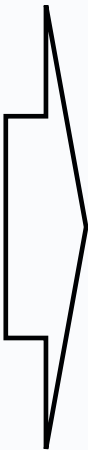
SSD >1TB (2020)

Datenbank

54		
48		
49		



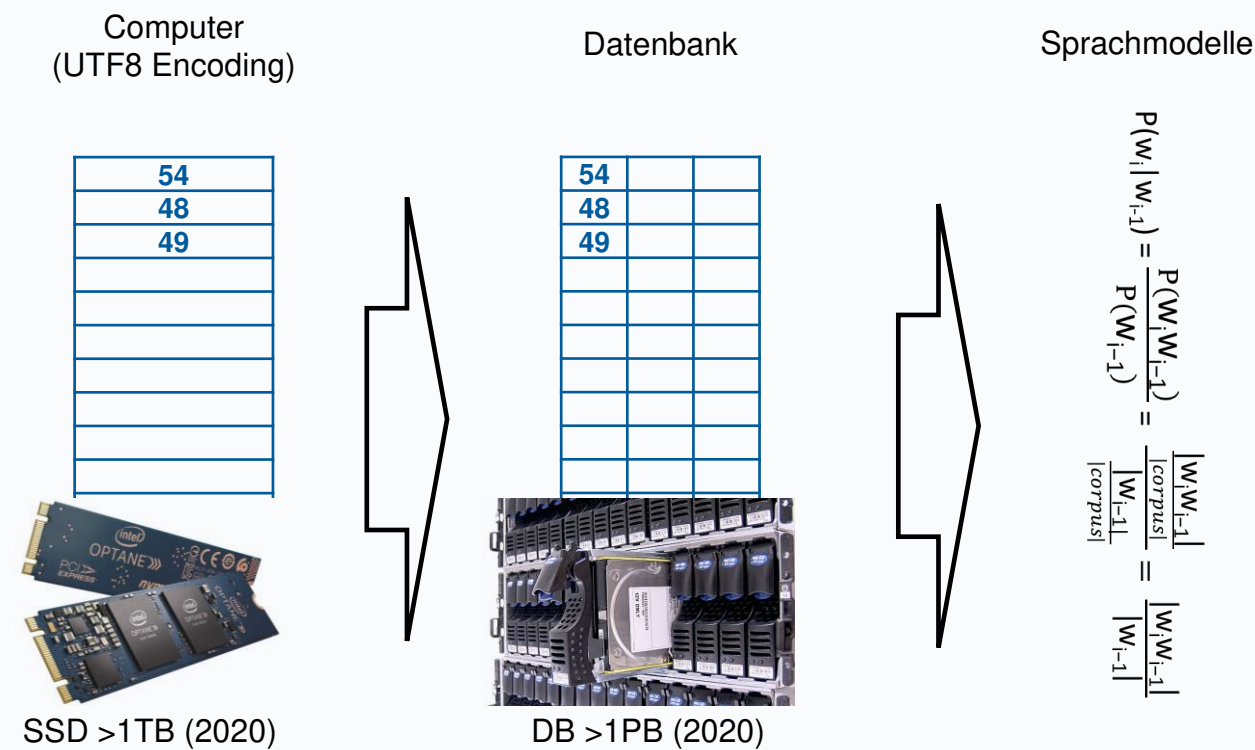
DB >1PB (2020)



# Word Embedding

From storing characters to storing relational data

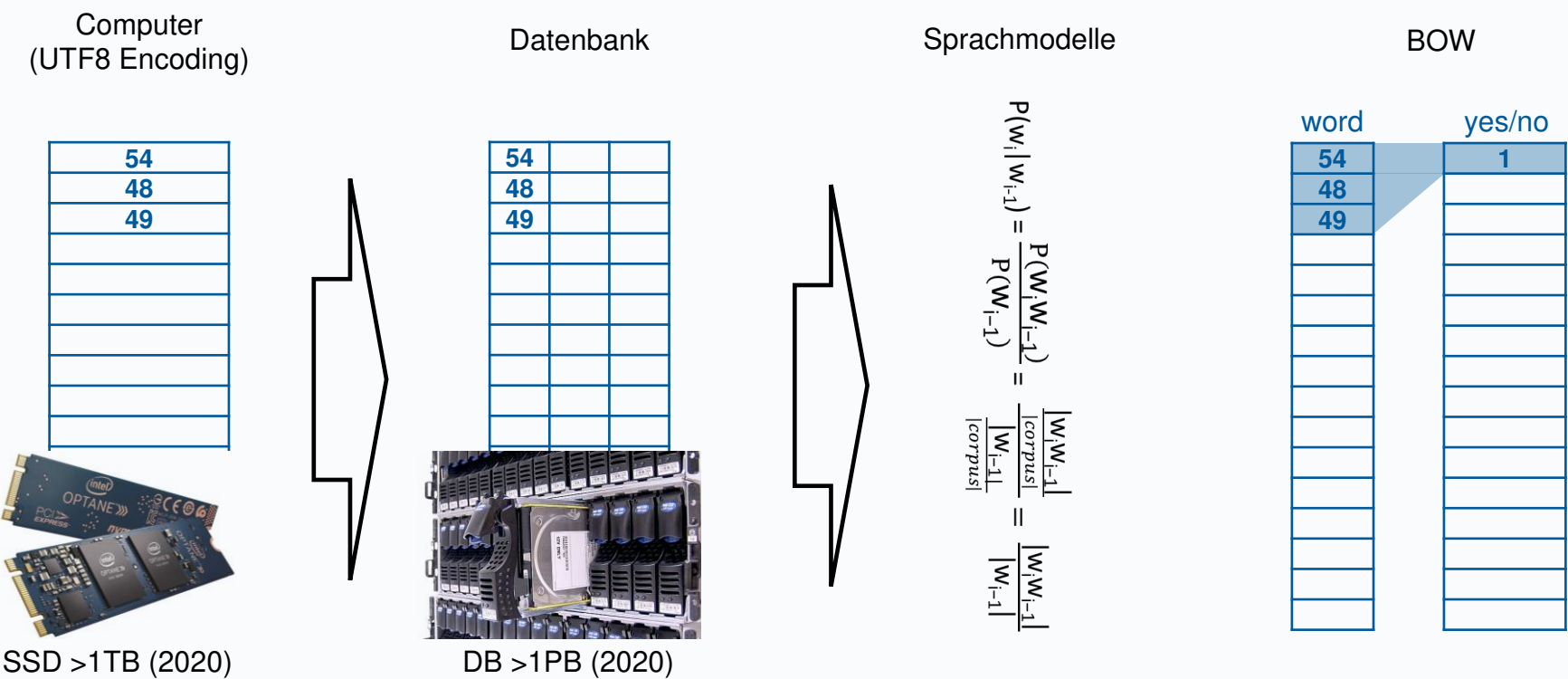
Data using a structure and language consistent with first-order predicate logic, first described in 1969 by Edgar F. Codd, where all data is represented in terms of tuples, grouped into relations.



# Word Embedding

From storing characters to storing relational data

Data using a structure and language consistent with first-order predicate logic, first described in 1969 by Edgar F. Codd, where all data is represented in terms of tuples, grouped into relations.



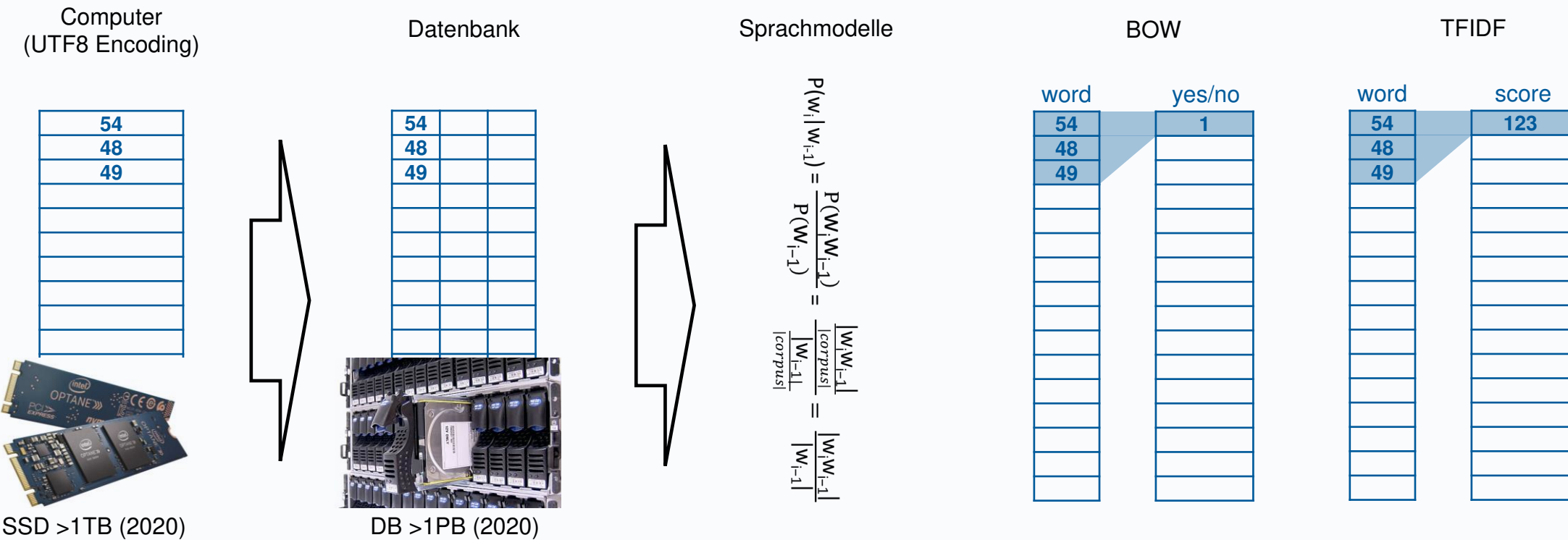


# Word Embedding

From storing characters to storing relational data



Data using a structure and language consistent with first-order predicate logic, first described in 1969 by Edgar F. Codd, where all data is represented in terms of tuples, grouped into relations.

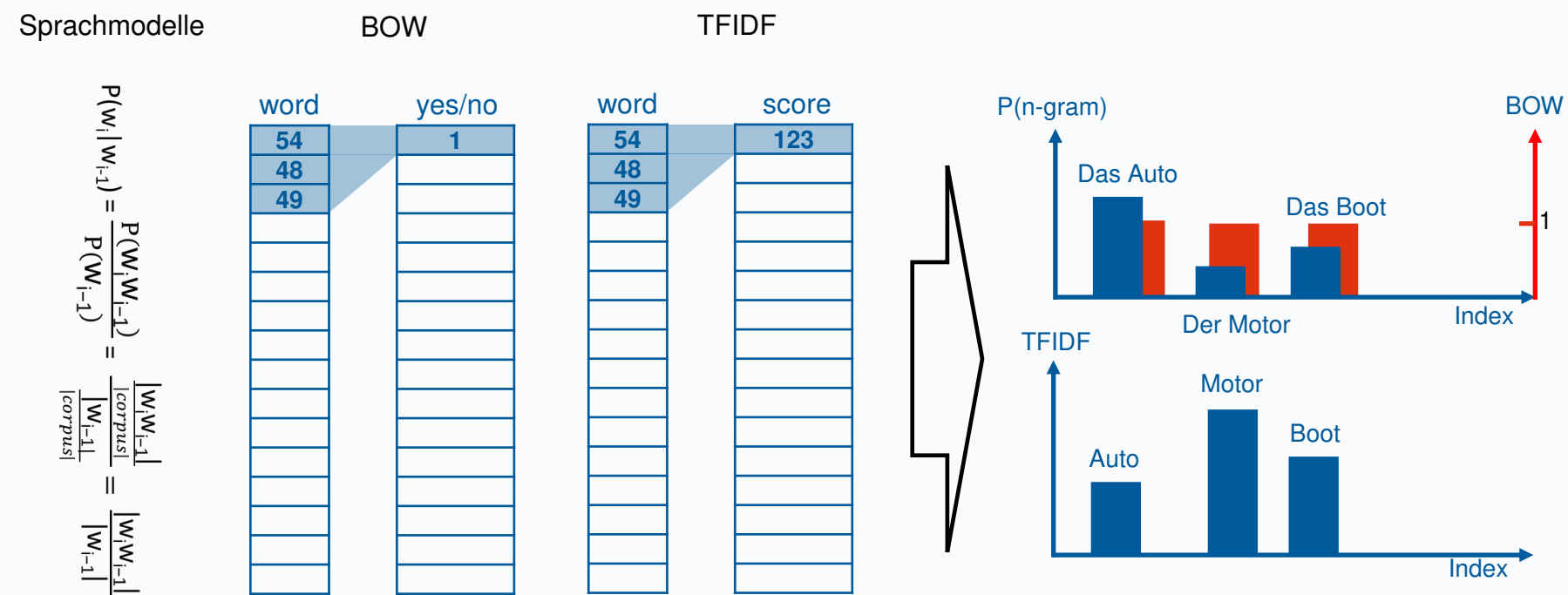




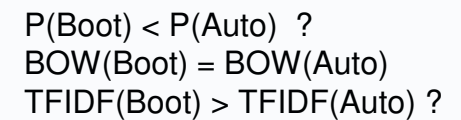
# Word Embedding

Mathematical view of relational data such as n-gram, BOW, TFIDF

Some information can be derived from n-gram, BOW, TFIDF key-value pairs.

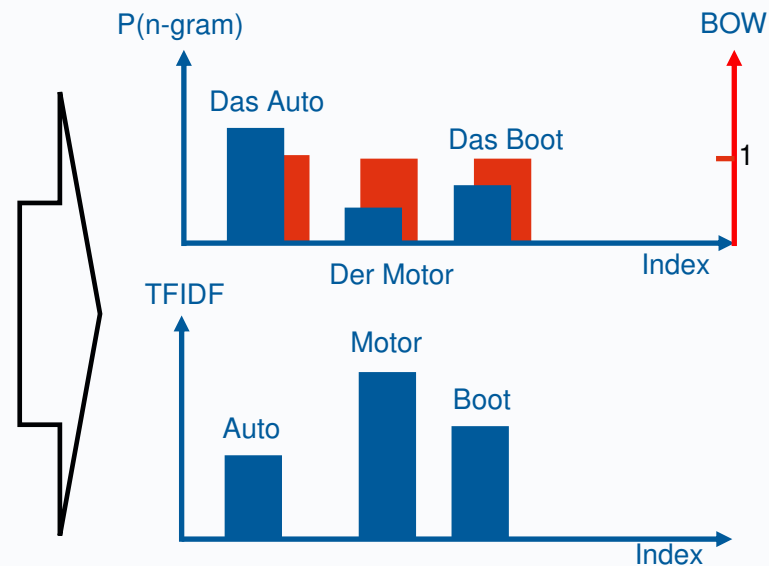


TFIDF

[illegible]

TFIDF

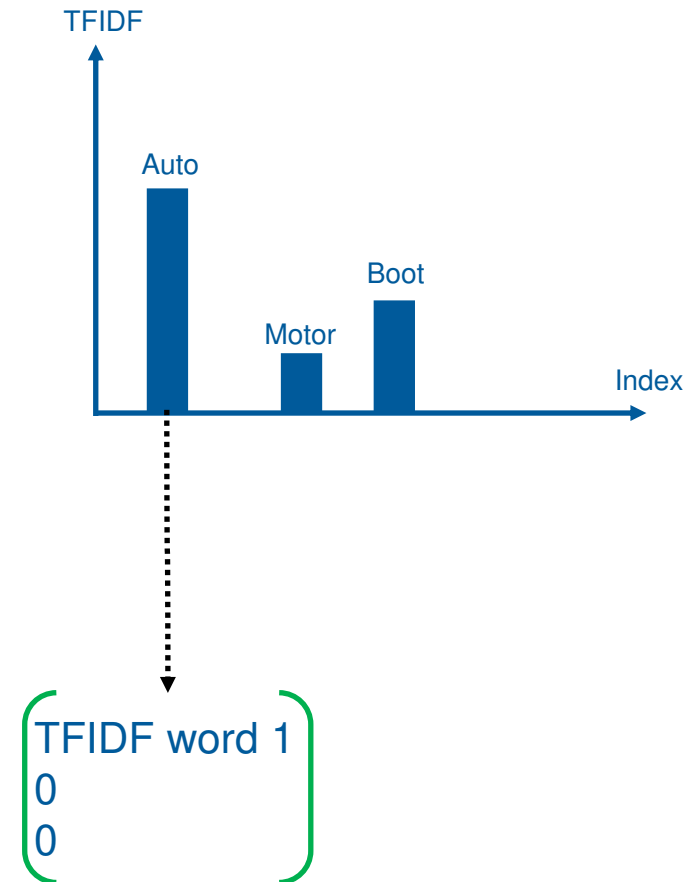
$$P(w_i | w_{i-1}) = \frac{P(w_i, w_{i-1})}{P(w_{i-1})} = \frac{\frac{|w_i, w_{i-1}|}{|corpus|}}{\frac{|w_{i-1}|}{|corpus|}} = \frac{|w_i, w_{i-1}|}{|w_{i-1}|}$$

[illegible][illegible]

P(Boot) < P(Auto) ?  
BOW(Boot) = BOW(Auto)  
TFIDF(Boot) > TFIDF(Auto) ?

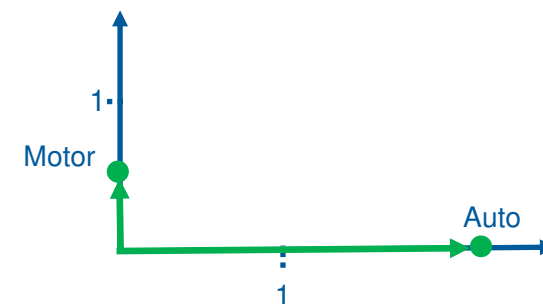
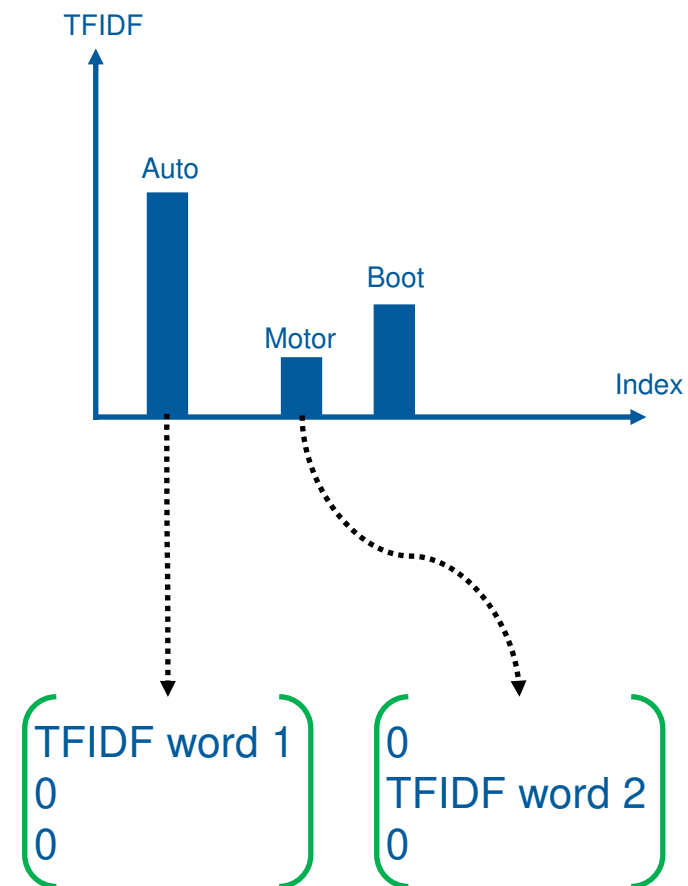
## Bestimmte Semantik ableitbar: Sprachmodell, Stop Words

Each word is represented by one unique vector.

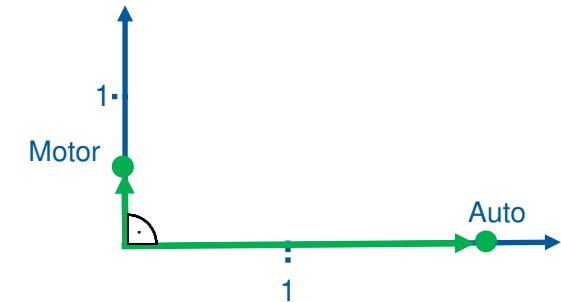
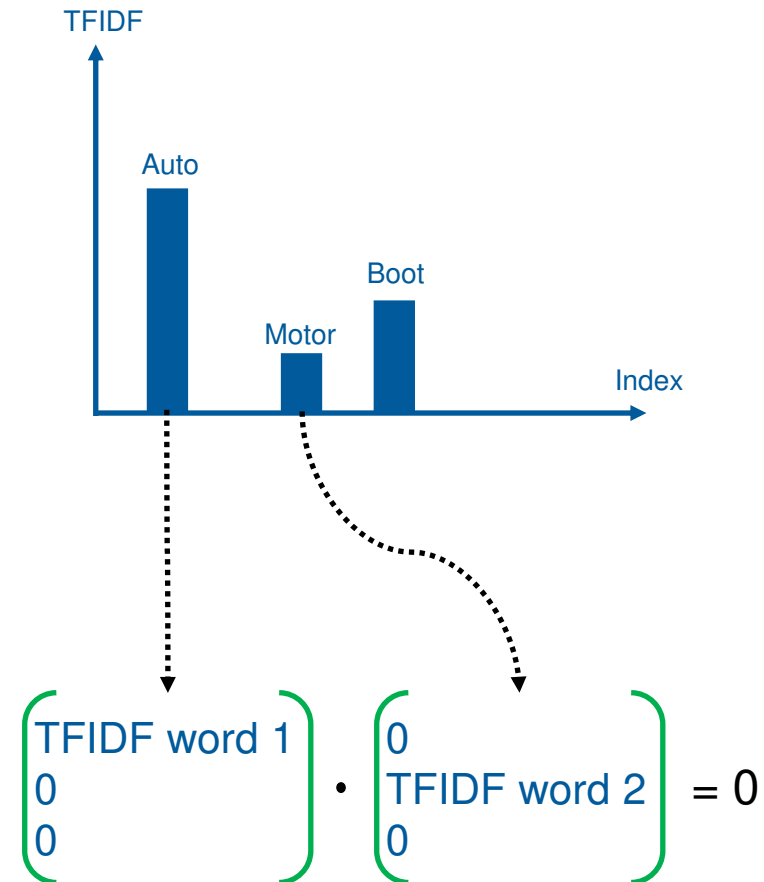




Each word is represented by one unique vector.



Defining word vectors in such a way that there is no semantic meaning between words.



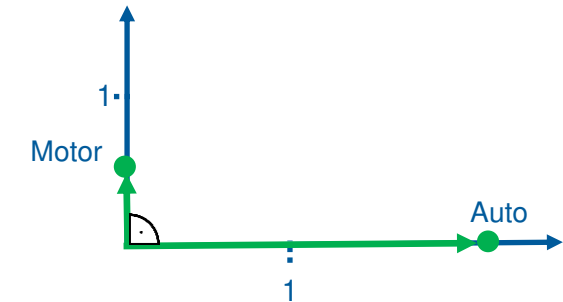
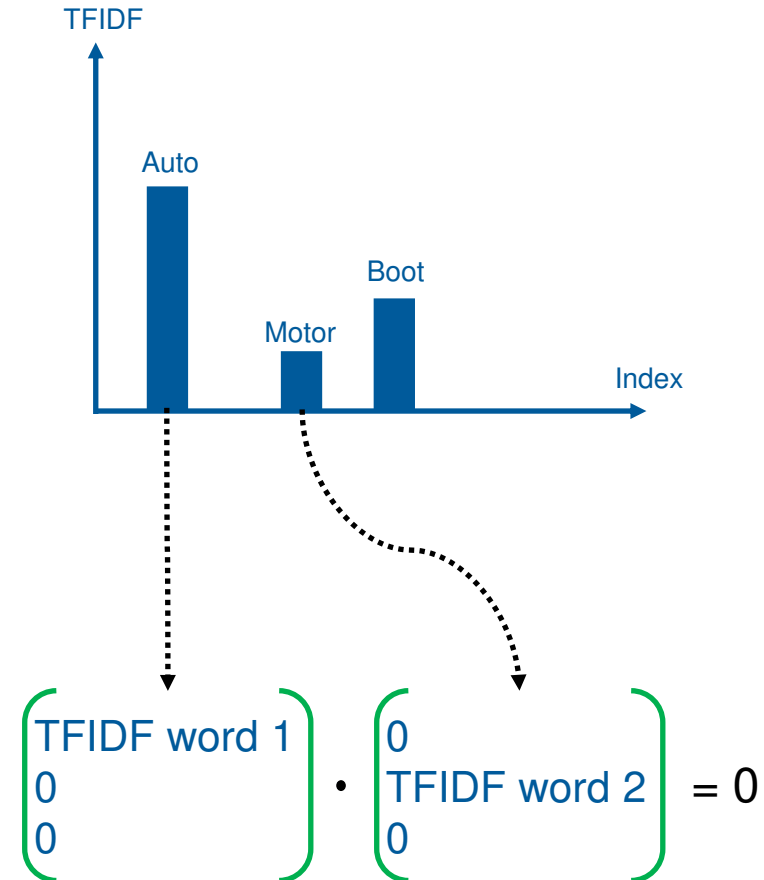
Dot-Product (zu Deut. Skalarprodukt):

$$\vec{a} \cdot \vec{b} = |\vec{a}| |\vec{b}| \cos \angle(\vec{a}, \vec{b})$$

$$\vec{a} \perp \vec{b} \iff \vec{a} \cdot \vec{b} = 0$$



Defining word vectors in such a way that there is no semantic meaning between words.



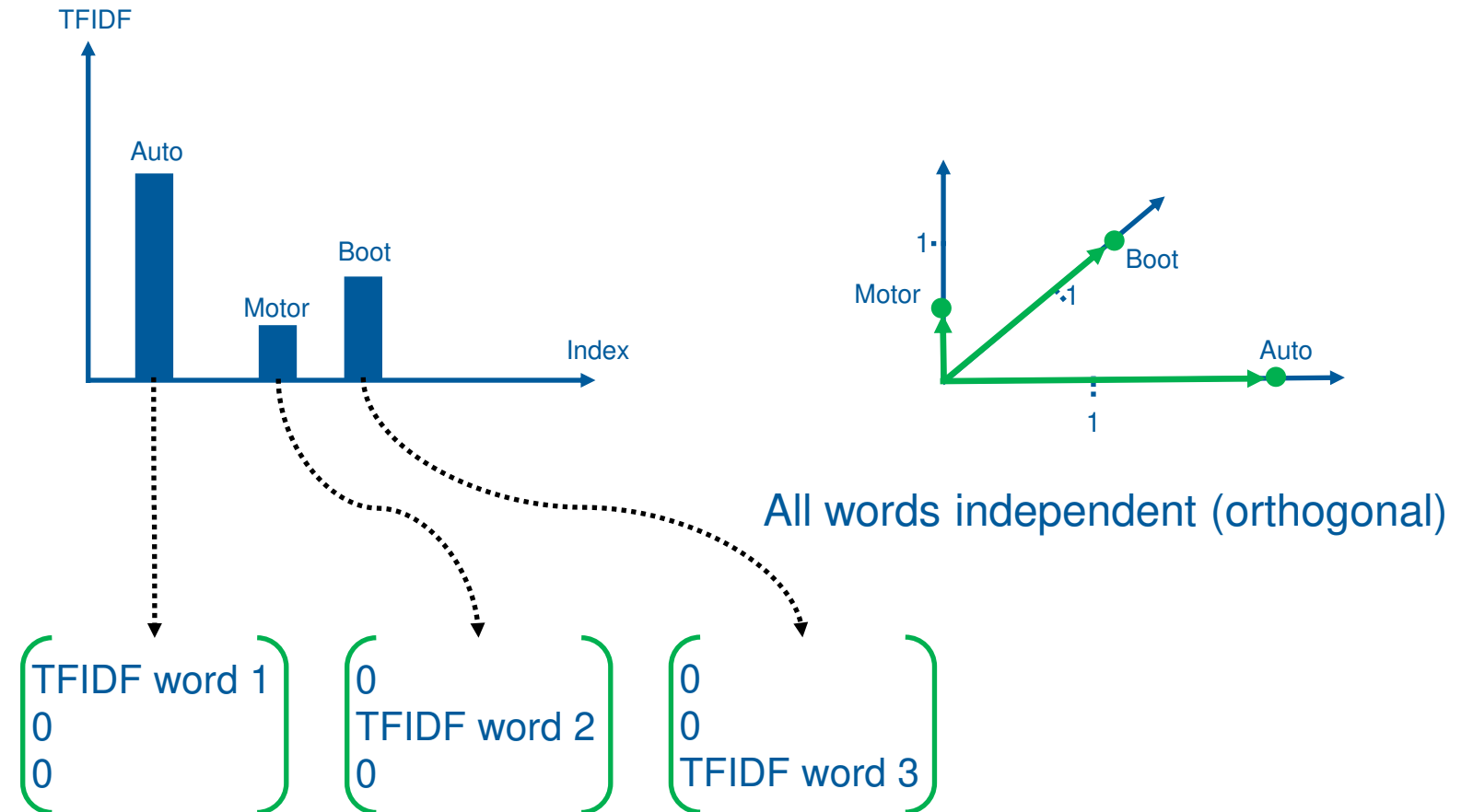
„Motor“ has nothing in common with „Auto“.

Dot-Product:

$$\vec{a} \cdot \vec{b} = |\vec{a}| |\vec{b}| \cos \angle(\vec{a}, \vec{b})$$

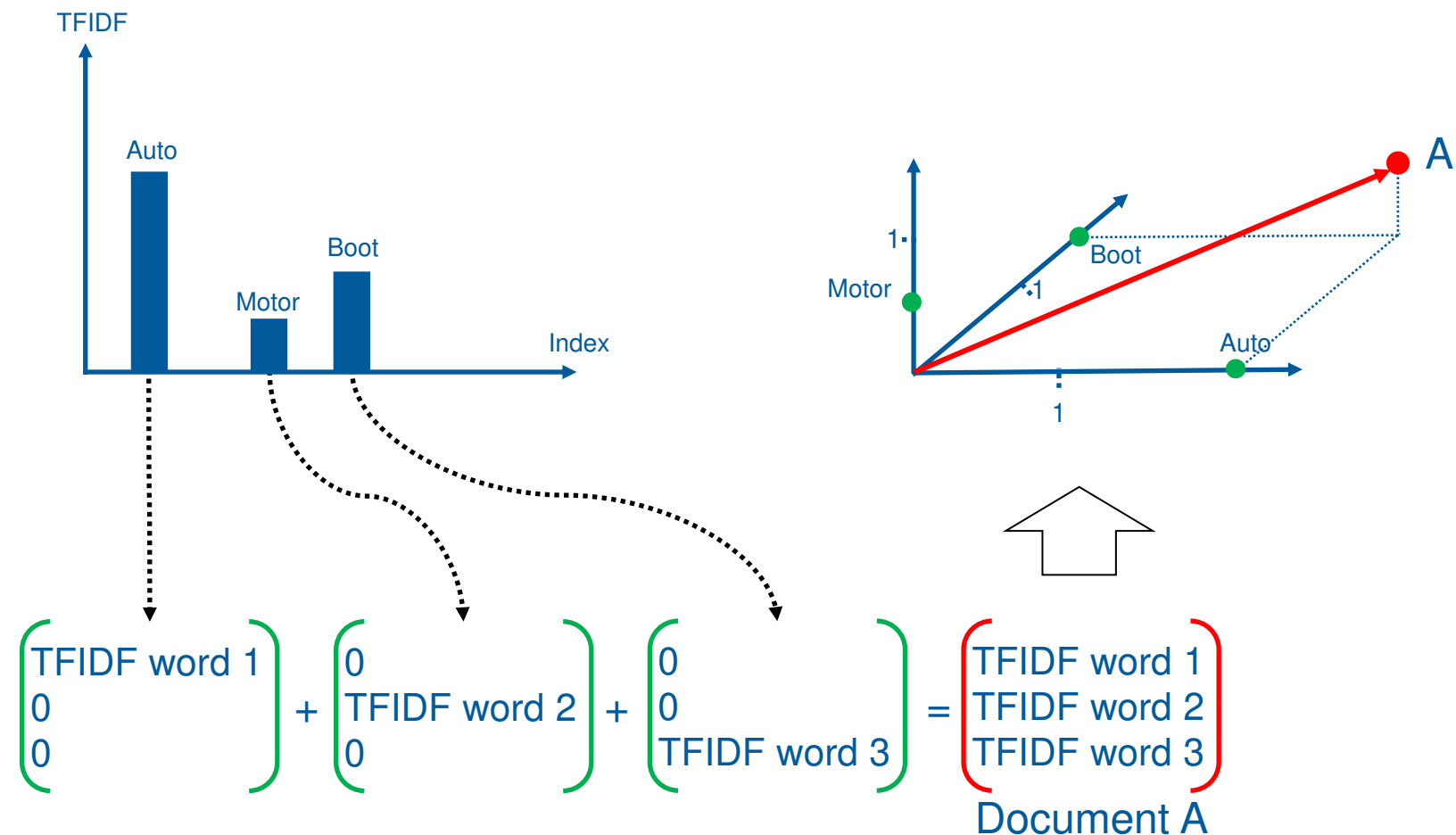
$$\vec{a} \perp \vec{b} \iff \vec{a} \cdot \vec{b} = 0$$

Defining word vectors in such a way that there is no semantic meaning between words.



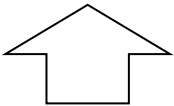
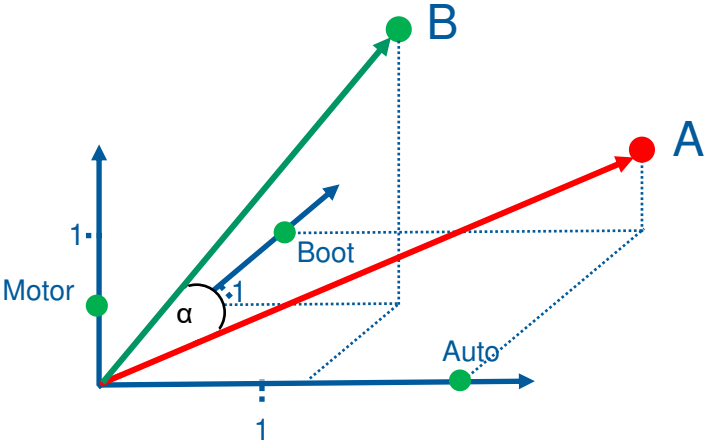
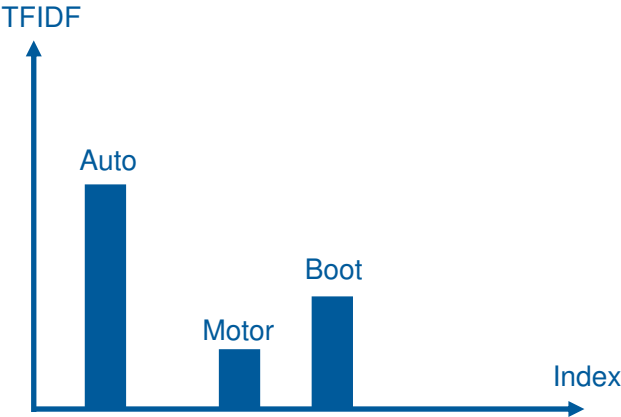


A document (collection of words) is the sum of all word vectors aka “document vector”





The similarity between documents is measured by the cosine similarity of the document vectors.



TFIDF word 1  
TFIDF word 2  
TFIDF word 3

Document B

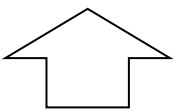
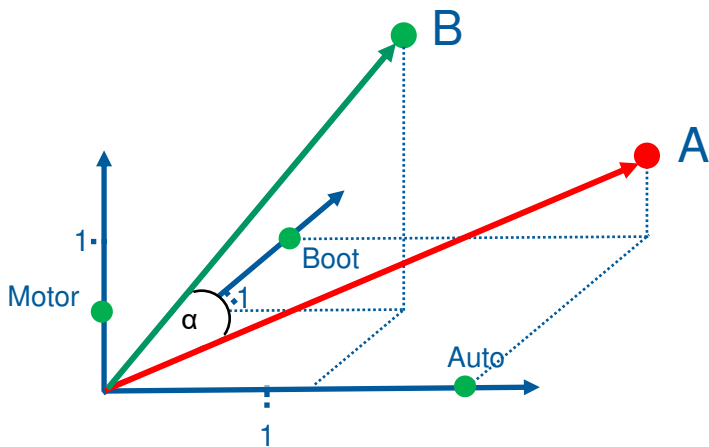
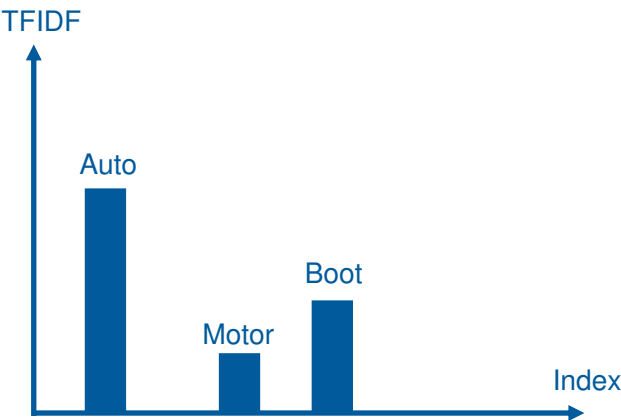
TFIDF word 1  
TFIDF word 2  
TFIDF word 3

Document A



The similarity between documents is measured by the cosine similarity of the document vectors.

We are using  
TFIDF and BOW  
for Document  
Classification,  
successfully.



TFIDF word 1  
TFIDF word 2  
TFIDF word 3

Document B

TFIDF word 1  
TFIDF word 2  
TFIDF word 3

Document A

# Word Embedding

## Words as Vectors



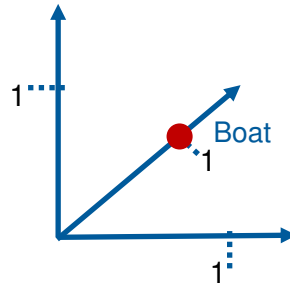
Word embedding is any of a set of language modeling and feature learning techniques in natural language processing where words from the vocabulary are mapped to vectors of real numbers.

$$\begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}$$

Boat

$$\begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}$$

Boat



# Word Embedding

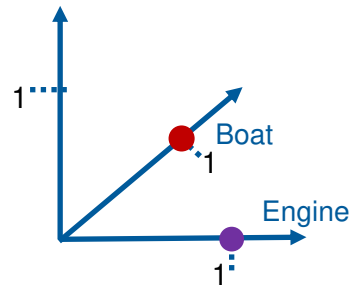
## Words as Vectors



Word embedding is any of a set of language modeling and feature learning techniques in natural language processing where words from the vocabulary are mapped to vectors of real numbers.

$$\begin{matrix} & \text{Engine} \\ \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} & - \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \\ \text{Boat} & \end{matrix}$$

$$\begin{matrix} & \text{Engine} \\ \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} & \cdot \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \\ \text{Boat} & \end{matrix}$$



# Word Embedding

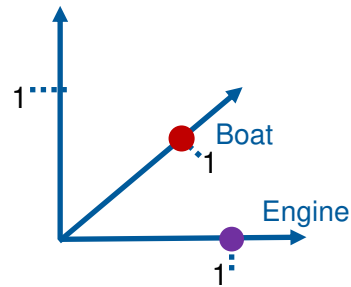
## Words as Vectors



Word embedding is any of a set of language modeling and feature learning techniques in natural language processing where words from the vocabulary are mapped to vectors of real numbers.

$$\begin{matrix} \text{Engine} \\ \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} - \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} = \begin{pmatrix} 1 \\ -1 \\ 0 \end{pmatrix} \\ \text{Boat} \end{matrix}$$

$$\begin{matrix} \text{Engine} \\ \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \cdot \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} = 0 \\ \text{Boat} \end{matrix}$$





# Word Embedding

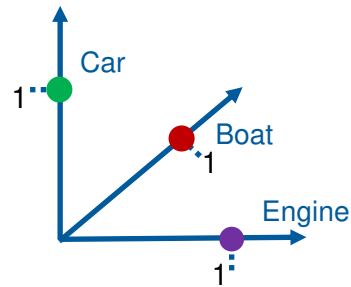
## Words as Vectors



Word embedding is any of a set of language modeling and feature learning techniques in natural language processing where words from the vocabulary are mapped to vectors of real numbers.

$$\begin{array}{c} \text{Engine} \\ \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} - \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} = \begin{pmatrix} 1 \\ -1 \\ 0 \end{pmatrix} \neq \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} \\ \text{Boat} \qquad \qquad \text{Car} \end{array}$$

$$\begin{array}{c} \text{Engine} \\ \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \cdot \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} \\ \text{Boat} \qquad \qquad \text{Car} \end{array}$$



# Word Embedding

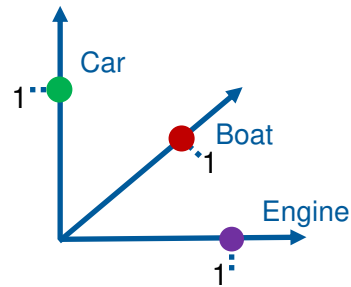
## Words as Vectors



Word embedding is any of a set of language modeling and feature learning techniques in natural language processing where words from the vocabulary are mapped to vectors of real numbers.

$$\begin{array}{c} \text{Engine} \\ \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} - \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} = \begin{pmatrix} 1 \\ -1 \\ 0 \end{pmatrix} \neq \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} - \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} = \begin{pmatrix} 0 \\ -1 \\ 1 \end{pmatrix} \\ \text{Boat} \quad \text{Car} \end{array}$$

$$\begin{array}{c} \text{Engine} \\ \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \cdot \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} \cdot \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} = 0 \\ \text{Boat} \quad \text{Car} \end{array}$$



# Word Embedding

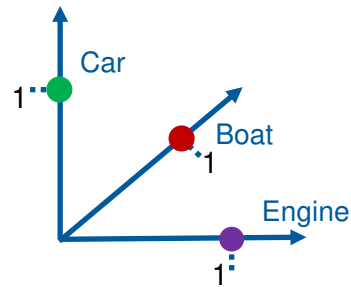
## Words as Vectors



Word embedding is any of a set of language modeling and feature learning techniques in natural language processing where words from the vocabulary are mapped to vectors of real numbers.

$$\begin{array}{c} \text{Engine} \\ \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} - \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} = \begin{pmatrix} 1 \\ -1 \\ 0 \end{pmatrix} \neq \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} - \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} = \begin{pmatrix} 0 \\ -1 \\ 1 \end{pmatrix} \\ \text{Boat} \quad \text{Car} \end{array}$$

$$\begin{array}{c} \text{Engine} \\ \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \cdot \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} \cdot \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \cdot \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} = 0 \\ \text{Boat} \quad \text{Car} \quad \text{Boat} \quad \text{Car} \end{array}$$



# Word Embedding

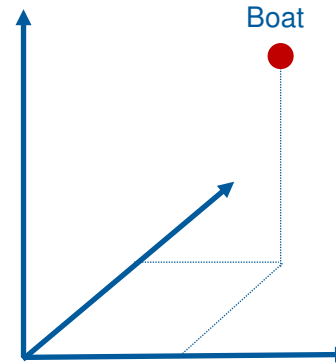
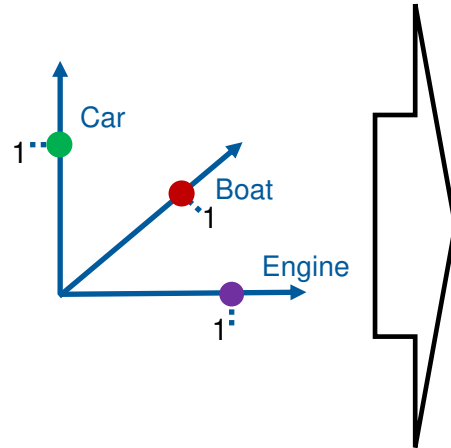
## Words as Vectors



Word embedding is any of a set of language modeling and feature learning techniques in natural language processing where words from the vocabulary are mapped to vectors of real numbers.

$$\begin{matrix} \text{Engine} \\ \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \end{matrix} - \begin{matrix} \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \\ \text{Boat} \end{matrix} = \begin{pmatrix} 1 \\ -1 \\ 0 \end{pmatrix} \neq \begin{matrix} \text{Engine} \\ \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} \end{matrix} - \begin{matrix} \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \\ \text{Car} \end{matrix} = \begin{pmatrix} 0 \\ -1 \\ 1 \end{pmatrix}$$

$$\begin{matrix} \text{Engine} \\ \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \end{matrix} \cdot \begin{matrix} \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \\ \text{Car} \end{matrix} = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} \cdot \begin{matrix} \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \\ \text{Car} \end{matrix} = \begin{matrix} \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \cdot \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} \\ \text{Boat} \end{matrix} = 0$$



$$\begin{pmatrix} 0.11 \\ 0.31 \\ 0.46 \end{pmatrix}$$

Boat

$$\begin{pmatrix} 0.11 \\ 0.31 \\ 0.46 \end{pmatrix}$$

Boat

# Word Embedding

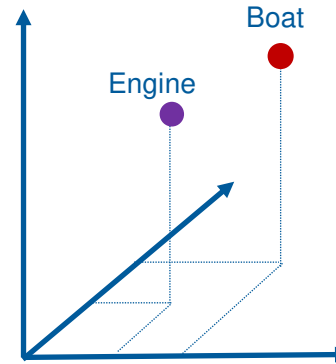
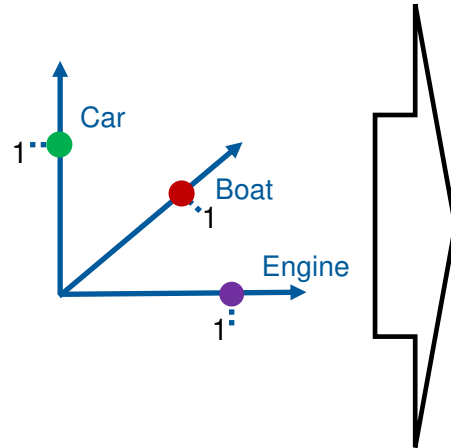
## Words as Vectors



Word embedding is any of a set of language modeling and feature learning techniques in natural language processing where words from the vocabulary are mapped to vectors of real numbers.

$$\begin{array}{c} \text{Engine} \\ \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} - \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} = \begin{pmatrix} 1 \\ -1 \\ 0 \end{pmatrix} \neq \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} - \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} = \begin{pmatrix} 0 \\ -1 \\ 1 \end{pmatrix} \\ \text{Boat} \quad \text{Car} \end{array}$$

$$\begin{array}{c} \text{Engine} \\ \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \cdot \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} \cdot \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \cdot \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} = 0 \\ \text{Boat} \quad \text{Car} \quad \text{Car} \end{array}$$



$$\begin{array}{c} \text{Engine} \\ \begin{pmatrix} 0.11 \\ 0.31 \\ 0.46 \end{pmatrix} - \begin{pmatrix} 0.05 \\ 0.01 \\ 0.21 \end{pmatrix} = \begin{pmatrix} 0.06 \\ 0.30 \\ 0.25 \end{pmatrix} \\ \text{Boat} \end{array}$$

$$\begin{array}{c} \text{Engine} \\ \begin{pmatrix} 0.11 \\ 0.31 \\ 0.46 \end{pmatrix} \cdot \begin{pmatrix} 0.43 \\ 0.01 \\ 0.21 \end{pmatrix} \\ \text{Boat} \end{array}$$

# Word Embedding

## Words as Vectors



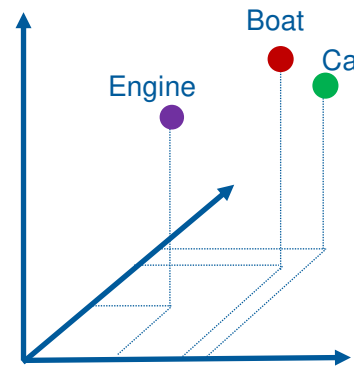
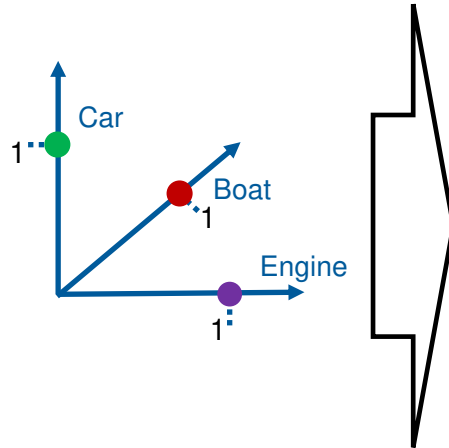
Word embedding is any of a set of language modeling and feature learning techniques in natural language processing where words from the vocabulary are mapped to vectors of real numbers.

$$\begin{array}{c} \text{Engine} \\ \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} - \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ -1 \\ 0 \end{bmatrix} \neq \begin{array}{c} \text{Engine} \\ \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} - \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ -1 \\ 1 \end{bmatrix} \end{array}$$

Boat Car

$$\begin{array}{c} \text{Engine} \\ \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \cdot \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = \begin{array}{c} \text{Engine} \\ \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \cdot \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = \begin{array}{c} \text{Boat} \\ \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \cdot \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} = 0 \end{array} \end{array}$$

Boat Car Car



$$\begin{array}{c} \text{Engine} \\ \begin{bmatrix} 0.11 \\ 0.31 \\ 0.46 \end{bmatrix} - \begin{bmatrix} 0.05 \\ 0.01 \\ 0.21 \end{bmatrix} = \begin{bmatrix} 0.06 \\ 0.30 \\ 0.25 \end{bmatrix} \end{array} \quad \begin{array}{c} \begin{bmatrix} 0.22 \\ 0.34 \\ 0.45 \end{bmatrix} \\ \text{Car} \end{array}$$

Boat

$$\begin{array}{c} \text{Engine} \\ \begin{bmatrix} 0.11 \\ 0.31 \\ 0.46 \end{bmatrix} \cdot \begin{bmatrix} 0.43 \\ 0.01 \\ 0.21 \end{bmatrix} = \begin{array}{c} \begin{bmatrix} 0.22 \\ 0.34 \\ 0.45 \end{bmatrix} \\ \text{Car} \end{array}$$

Boat

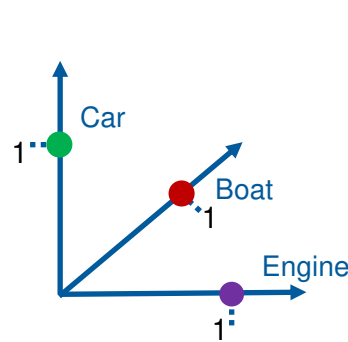
# Word Embedding

## Words as Vectors

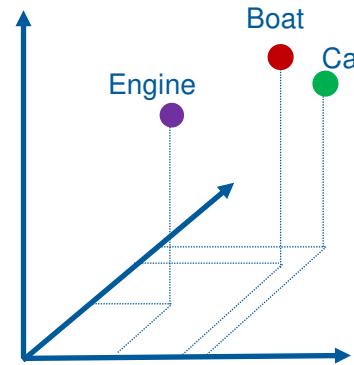


Word embedding is any of a set of language modeling and feature learning techniques in natural language processing where words from the vocabulary are mapped to vectors of real numbers.

$$\begin{matrix} \text{Engine} \\ \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \end{matrix} - \begin{matrix} \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \\ \text{Boat} \end{matrix} = \begin{bmatrix} 1 \\ -1 \\ 0 \end{bmatrix} \neq \begin{matrix} \text{Engine} \\ \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \end{matrix} - \begin{matrix} \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \\ \text{Car} \end{matrix} = \begin{bmatrix} 0 \\ -1 \\ 1 \end{bmatrix}$$



$$\begin{matrix} \text{Engine} \\ \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \end{matrix} \cdot \begin{matrix} \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \\ \text{Car} \end{matrix} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \cdot \begin{matrix} \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \\ \text{Car} \end{matrix} = \begin{matrix} \text{Boat} \\ \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \end{matrix} \cdot \begin{matrix} \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \\ \text{Car} \end{matrix} = 0$$

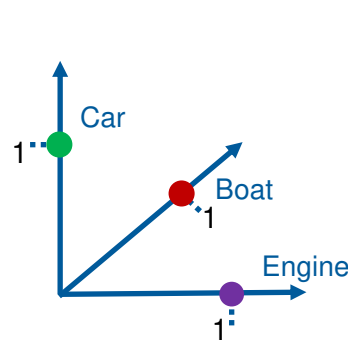


$$\begin{matrix} \text{Engine} \\ \begin{bmatrix} 0.11 \\ 0.31 \\ 0.46 \end{bmatrix} \end{matrix} - \begin{matrix} \begin{bmatrix} 0.05 \\ 0.01 \\ 0.21 \end{bmatrix} \\ \text{Boat} \end{matrix} = \begin{bmatrix} 0.06 \\ 0.30 \\ 0.25 \end{bmatrix} \neq \begin{matrix} \text{Engine} \\ \begin{bmatrix} 0.22 \\ 0.34 \\ 0.45 \end{bmatrix} \end{matrix} - \begin{matrix} \begin{bmatrix} 0.05 \\ 0.01 \\ 0.21 \end{bmatrix} \\ \text{Car} \end{matrix} = \begin{bmatrix} 0.17 \\ 0.33 \\ 0.24 \end{bmatrix}$$

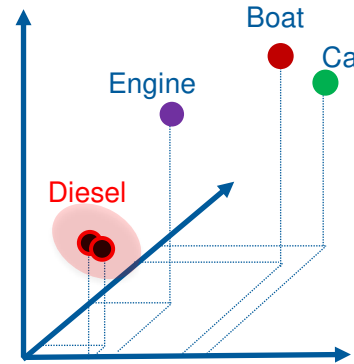
$$\begin{matrix} \text{Engine} \\ \begin{bmatrix} 0.11 \\ 0.31 \\ 0.46 \end{bmatrix} \end{matrix} \cdot \begin{matrix} \begin{bmatrix} 0.43 \\ 0.01 \\ 0.21 \end{bmatrix} \\ \text{Car} \end{matrix} \neq \begin{matrix} \text{Engine} \\ \begin{bmatrix} 0.22 \\ 0.34 \\ 0.45 \end{bmatrix} \end{matrix} \cdot \begin{matrix} \begin{bmatrix} 0.43 \\ 0.01 \\ 0.21 \end{bmatrix} \\ \text{Car} \end{matrix} \neq \begin{matrix} \text{Boat} \\ \begin{bmatrix} 0.11 \\ 0.31 \\ 0.46 \end{bmatrix} \end{matrix} \cdot \begin{matrix} \begin{bmatrix} 0.22 \\ 0.34 \\ 0.45 \end{bmatrix} \\ \text{Car} \end{matrix} \neq 0$$

Word embedding is any of a set of language modeling and feature learning techniques in natural language processing where words from the vocabulary are mapped to vectors of real numbers.

$$\begin{matrix} \text{Engine} \\ \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \\ \text{Boat} \end{matrix} - \begin{matrix} \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \\ \text{Car} \end{matrix} = \begin{bmatrix} 1 \\ -1 \\ 0 \end{bmatrix} \neq \begin{matrix} \text{Engine} \\ \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \\ \text{Car} \end{matrix} - \begin{matrix} \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \\ \text{Car} \end{matrix} = \begin{bmatrix} 0 \\ -1 \\ 1 \end{bmatrix}$$



$$\begin{matrix} \text{Engine} \\ \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \\ \text{Boat} \end{matrix} \cdot \begin{matrix} \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \\ \text{Car} \end{matrix} = \begin{matrix} \text{Engine} \\ \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \\ \text{Car} \end{matrix} \cdot \begin{matrix} \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \\ \text{Car} \end{matrix} = \begin{matrix} \text{Boat} \\ \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \\ \text{Car} \end{matrix} \cdot \begin{matrix} \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \\ \text{Car} \end{matrix} = 0$$



$$\begin{matrix} \text{Engine} \\ \begin{bmatrix} 0.11 \\ 0.31 \\ 0.46 \end{bmatrix} \\ \text{Boat} \end{matrix} - \begin{matrix} \begin{bmatrix} 0.05 \\ 0.01 \\ 0.21 \end{bmatrix} \\ \text{Car} \end{matrix} = \begin{matrix} \text{Close to Diesel} \\ \begin{bmatrix} 0.06 \\ 0.30 \\ 0.25 \end{bmatrix} \end{matrix} \neq \begin{matrix} \text{Engine} \\ \begin{bmatrix} 0.22 \\ 0.34 \\ 0.45 \end{bmatrix} \\ \text{Car} \end{matrix} - \begin{matrix} \begin{bmatrix} 0.05 \\ 0.01 \\ 0.21 \end{bmatrix} \\ \text{Car} \end{matrix} = \begin{matrix} \text{Close to Diesel} \\ \begin{bmatrix} 0.17 \\ 0.33 \\ 0.24 \end{bmatrix} \end{matrix}$$

$$\begin{matrix} \text{Engine} \\ \begin{bmatrix} 0.11 \\ 0.31 \\ 0.46 \end{bmatrix} \\ \text{Boat} \end{matrix} \cdot \begin{matrix} \begin{bmatrix} 0.43 \\ 0.01 \\ 0.21 \end{bmatrix} \\ \text{Car} \end{matrix} \neq \begin{matrix} \text{Engine} \\ \begin{bmatrix} 0.22 \\ 0.34 \\ 0.45 \end{bmatrix} \\ \text{Car} \end{matrix} \cdot \begin{matrix} \begin{bmatrix} 0.43 \\ 0.01 \\ 0.21 \end{bmatrix} \\ \text{Car} \end{matrix} \neq \begin{matrix} \text{Boat} \\ \begin{bmatrix} 0.11 \\ 0.31 \\ 0.46 \end{bmatrix} \\ \text{Car} \end{matrix} \cdot \begin{matrix} \begin{bmatrix} 0.22 \\ 0.34 \\ 0.45 \end{bmatrix} \\ \text{Car} \end{matrix} \neq 0$$



# Word Embedding

## Words as Vectors



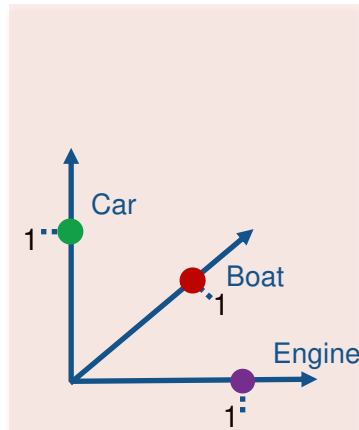
Word embedding is any of a set of language modeling and feature learning techniques in natural language processing where words from the vocabulary are mapped to vectors of real numbers.

$$\begin{array}{c} \text{Engine} \\ \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} - \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ -1 \\ 0 \end{bmatrix} \neq \begin{array}{c} \text{Engine} \\ \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} - \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ -1 \\ 1 \end{bmatrix} \end{array}$$

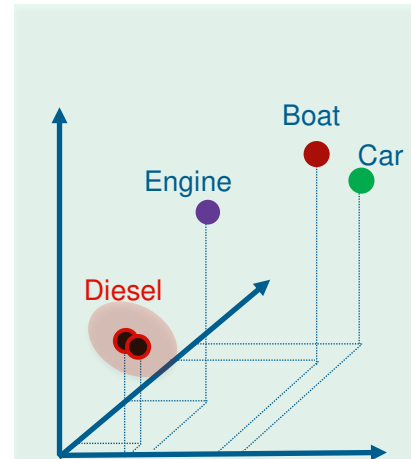
Boat Car

$$\begin{array}{c} \text{Engine} \\ \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \cdot \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \cdot \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = \begin{array}{c} \text{Boat} \\ \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \cdot \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} = 0 \end{array}$$

Boat Car Car



at most 3 words



maximal number of words?

$$\begin{array}{c} \text{Engine} \\ \begin{bmatrix} 0.11 \\ 0.31 \\ 0.46 \end{bmatrix} - \begin{bmatrix} 0.05 \\ 0.01 \\ 0.21 \end{bmatrix} = \begin{array}{c} \text{Close to Diesel} \\ \begin{bmatrix} 0.06 \\ 0.30 \\ 0.25 \end{bmatrix} \end{array} \neq \begin{array}{c} \text{Engine} \\ \begin{bmatrix} 0.22 \\ 0.34 \\ 0.45 \end{bmatrix} - \begin{bmatrix} 0.05 \\ 0.01 \\ 0.21 \end{bmatrix} = \begin{array}{c} \text{Close to Diesel} \\ \begin{bmatrix} 0.17 \\ 0.33 \\ 0.24 \end{bmatrix} \end{array}$$

Boat Car

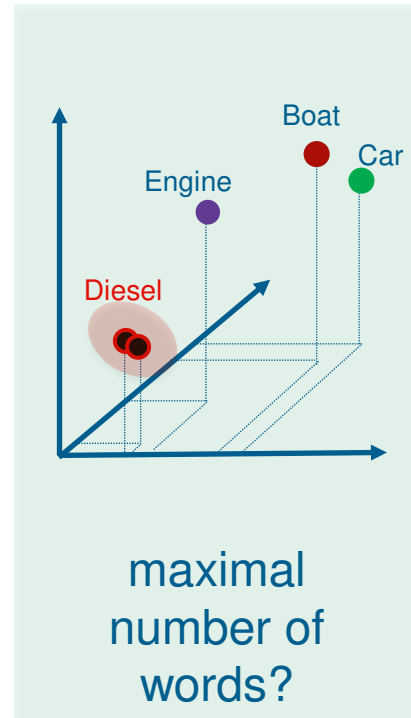
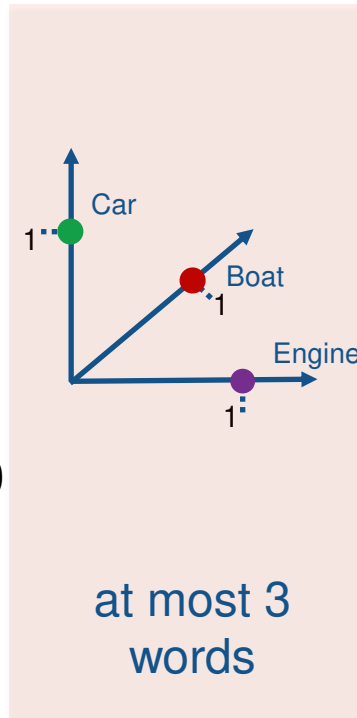
$$\begin{array}{c} \text{Engine} \\ \begin{bmatrix} 0.11 \\ 0.31 \\ 0.46 \end{bmatrix} \cdot \begin{bmatrix} 0.43 \\ 0.01 \\ 0.21 \end{bmatrix} \neq \begin{array}{c} \text{Engine} \\ \begin{bmatrix} 0.22 \\ 0.34 \\ 0.45 \end{bmatrix} \cdot \begin{bmatrix} 0.43 \\ 0.01 \\ 0.21 \end{bmatrix} \neq \begin{array}{c} \text{Boat} \\ \begin{bmatrix} 0.11 \\ 0.31 \\ 0.46 \end{bmatrix} \cdot \begin{bmatrix} 0.22 \\ 0.34 \\ 0.45 \end{bmatrix} \neq 0 \end{array}$$

Boat Car Car

Word embedding is any of a set of language modeling and feature learning techniques in natural language processing where words from the vocabulary are mapped to vectors of real numbers.

$$\begin{matrix} \text{Engine} \\ \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \end{matrix} - \begin{matrix} \text{Boat} \\ \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \end{matrix} = \begin{bmatrix} 1 \\ -1 \\ 0 \end{bmatrix} \neq \begin{matrix} \text{Engine} \\ \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \end{matrix} - \begin{matrix} \text{Car} \\ \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \end{matrix} = \begin{bmatrix} 0 \\ -1 \\ 1 \end{bmatrix}$$

$$\begin{matrix} \text{Engine} \\ \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \end{matrix} \cdot \begin{matrix} \text{Boat} \\ \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \end{matrix} = \begin{matrix} \text{Engine} \\ \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \end{matrix} \cdot \begin{matrix} \text{Car} \\ \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \end{matrix} = \begin{matrix} \text{Boat} \\ \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \end{matrix} \cdot \begin{matrix} \text{Car} \\ \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \end{matrix} = 0$$



$$\begin{matrix} \text{Boat} \\ \begin{bmatrix} 0.11 \\ 0.31 \\ 0.46 \end{bmatrix} \end{matrix} - \begin{matrix} \text{Engine} \\ \begin{bmatrix} 0.05 \\ 0.01 \\ 0.21 \end{bmatrix} \end{matrix} = \begin{matrix} \text{Close to Diesel} \\ \begin{bmatrix} 0.06 \\ 0.30 \\ 0.25 \end{bmatrix} \end{matrix} \neq \begin{matrix} \text{Car} \\ \begin{bmatrix} 0.22 \\ 0.34 \\ 0.45 \end{bmatrix} \end{matrix} - \begin{matrix} \text{Engine} \\ \begin{bmatrix} 0.05 \\ 0.01 \\ 0.21 \end{bmatrix} \end{matrix} = \begin{matrix} \text{Close to Diesel} \\ \begin{bmatrix} 0.17 \\ 0.33 \\ 0.24 \end{bmatrix} \end{matrix}$$

$$\begin{matrix} \text{Boat} \\ \begin{bmatrix} 0.11 \\ 0.31 \\ 0.46 \end{bmatrix} \end{matrix} \cdot \begin{matrix} \text{Engine} \\ \begin{bmatrix} 0.43 \\ 0.01 \\ 0.21 \end{bmatrix} \end{matrix} \neq \begin{matrix} \text{Car} \\ \begin{bmatrix} 0.22 \\ 0.34 \\ 0.45 \end{bmatrix} \end{matrix} \cdot \begin{matrix} \text{Engine} \\ \begin{bmatrix} 0.43 \\ 0.01 \\ 0.21 \end{bmatrix} \end{matrix} \neq \begin{matrix} \text{Boat} \\ \begin{bmatrix} 0.11 \\ 0.31 \\ 0.46 \end{bmatrix} \end{matrix} \cdot \begin{matrix} \text{Car} \\ \begin{bmatrix} 0.22 \\ 0.34 \\ 0.45 \end{bmatrix} \end{matrix} \neq 0$$

# Word Embedding

## Transform a Word Vector



Word embedding is any of a set of language modeling and feature learning techniques in natural language processing where words from the vocabulary are mapped to vectors of real numbers.

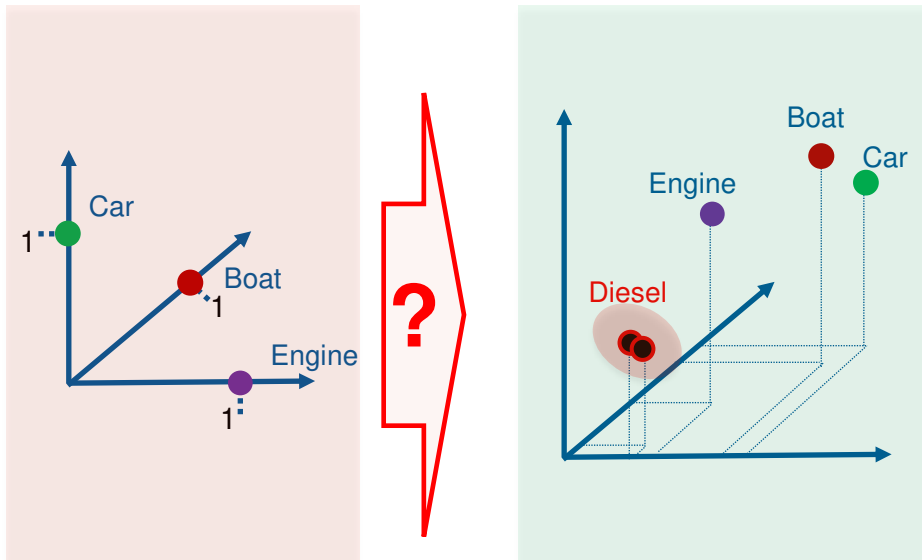


# Word Embedding

## Transform a Word Vector



Word embedding is any of a set of language modeling and feature learning techniques in natural language processing where words from the vocabulary are mapped to vectors of real numbers.

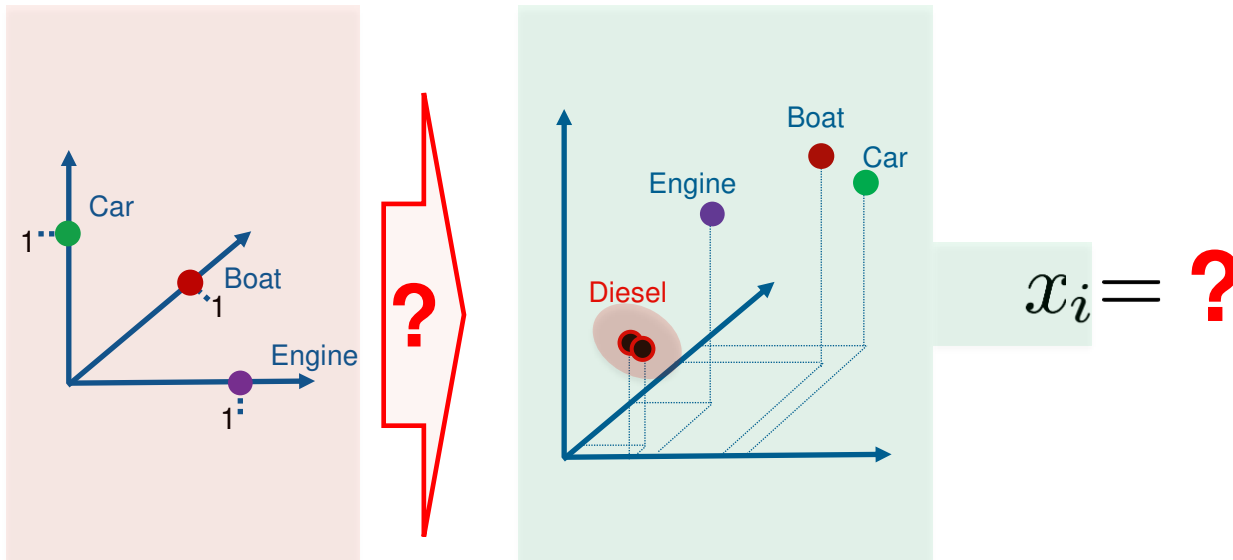


# Word Embedding

## Transform a Word Vector



Word embedding is any of a set of language modeling and feature learning techniques in natural language processing where words from the vocabulary are mapped to vectors of real numbers.

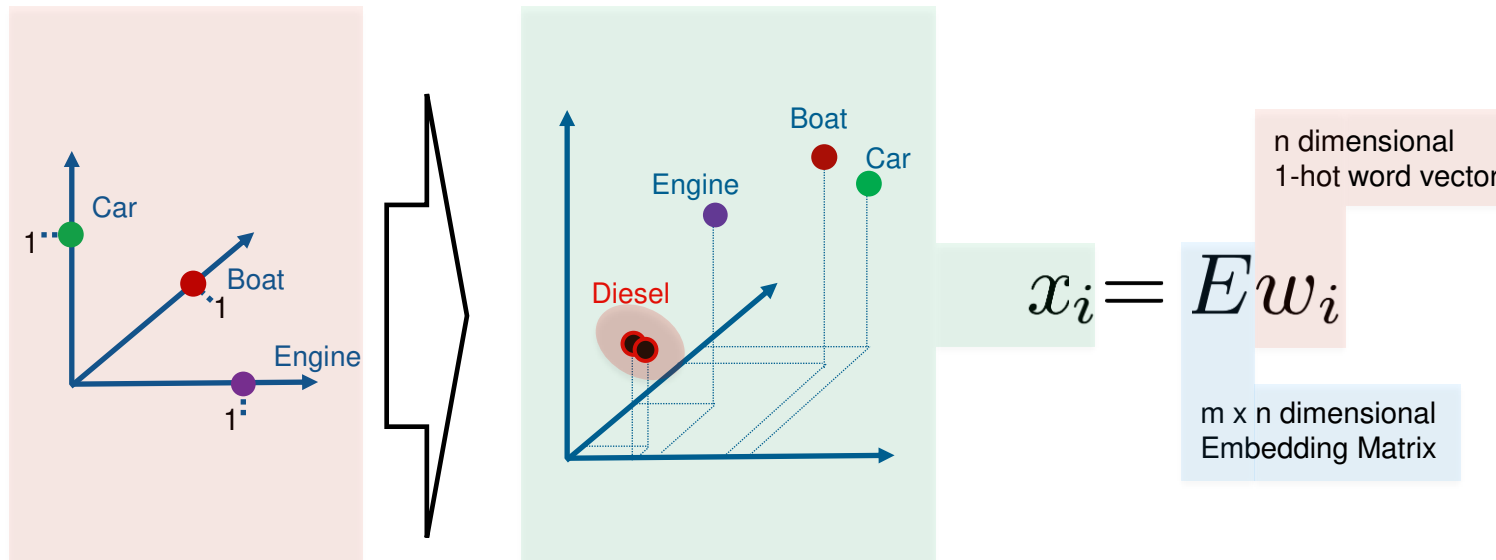


# Word Embedding

## Transform a Word Vector



Word embedding is any of a set of language modeling and feature learning techniques in natural language processing where words from the vocabulary are mapped to vectors of real numbers.





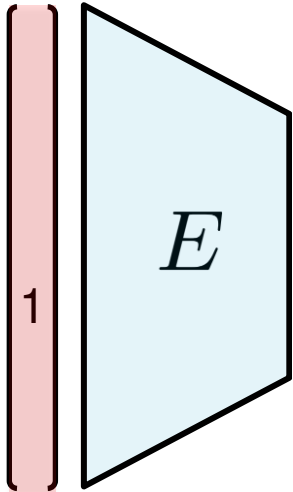
Word embedding is any of a set of language modeling and feature learning techniques in natural language processing where words from the vocabulary are mapped to vectors of real numbers.

$w_i$

1

Word embedding is any of a set of language modeling and feature learning techniques in natural language processing where words from the vocabulary are mapped to vectors of real numbers.

$$E w_i$$





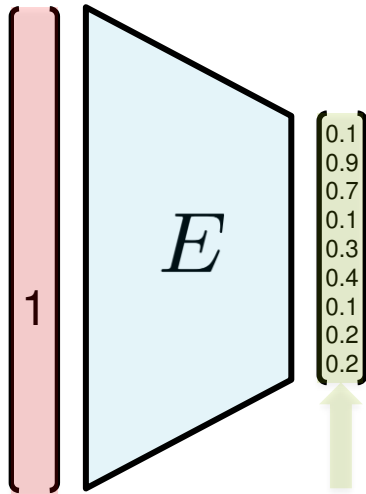
# Word Embedding

## Transform a Word Vector



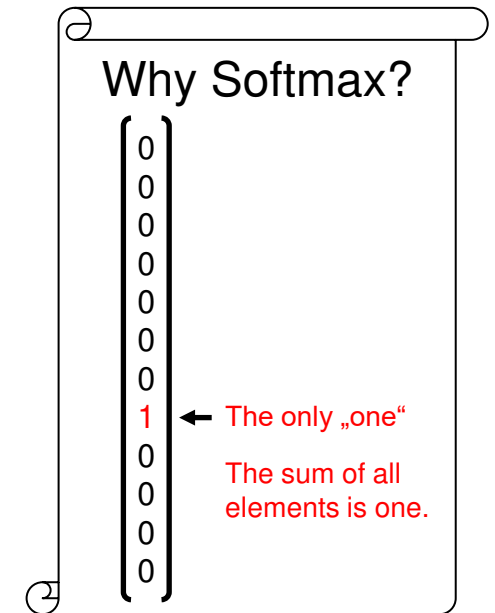
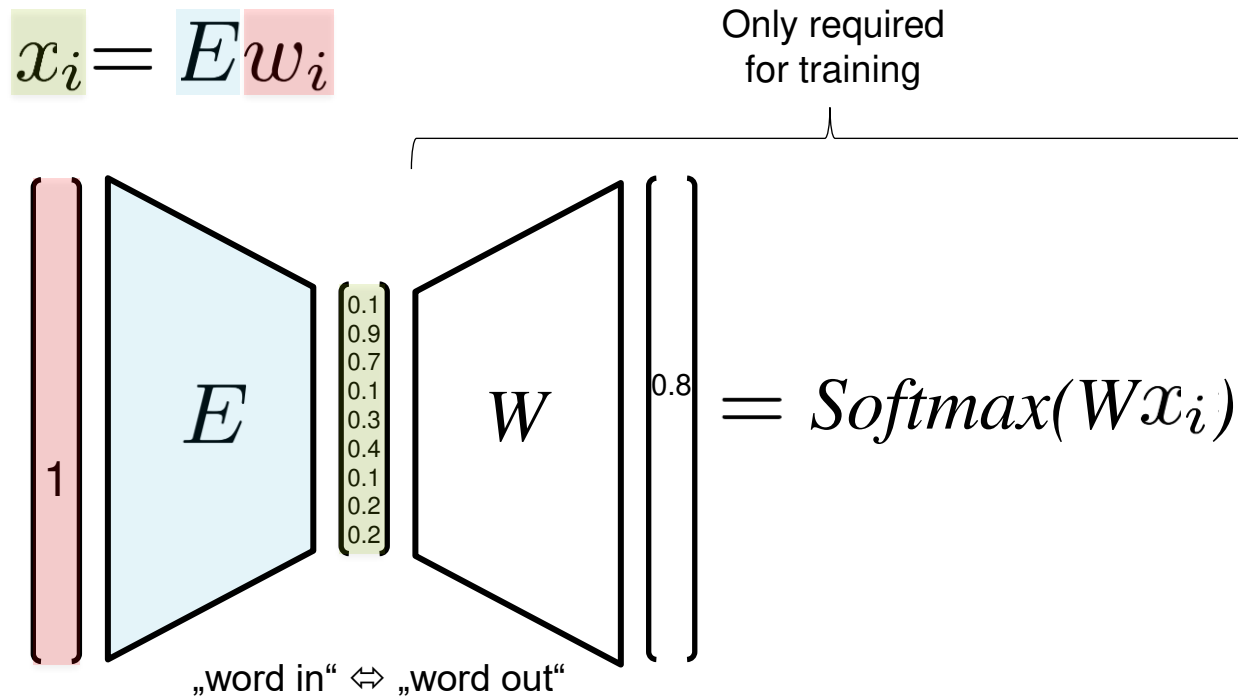
Word embedding is any of a set of language modeling and feature learning techniques in natural language processing where words from the vocabulary are mapped to vectors of real numbers.

$$x_i = E w_i$$



Word Embedding

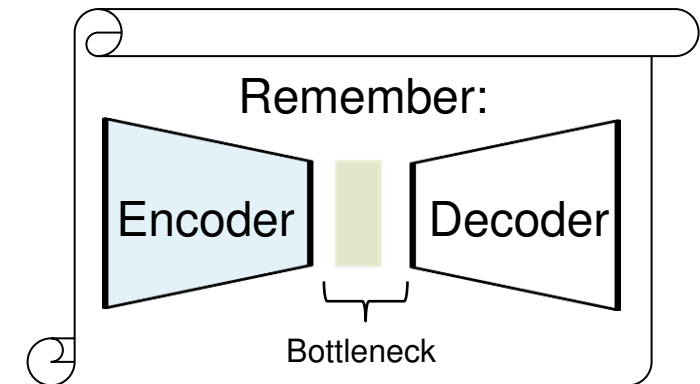
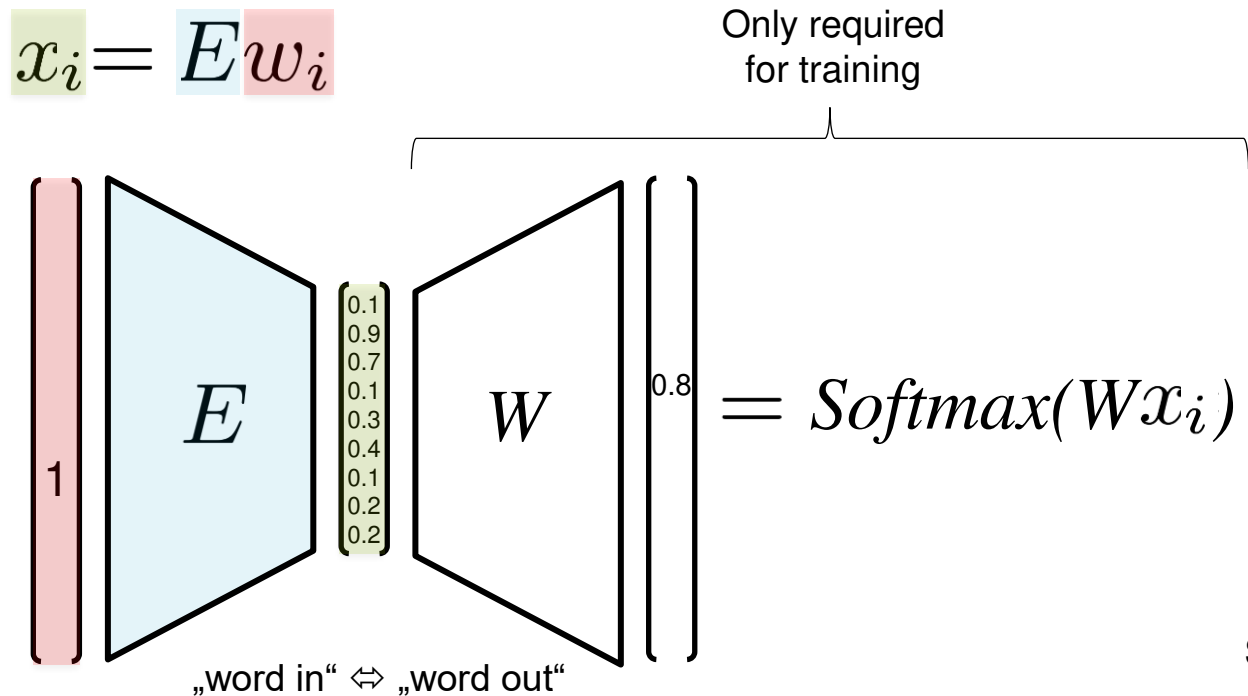
Word embedding is any of a set of language modeling and feature learning techniques in natural language processing where words from the vocabulary are mapped to vectors of real numbers.



Softmax:

$$\sigma(\mathbf{z})_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}} \text{ for } i = 1, \dots, K \text{ and } \mathbf{z} = (z_1, \dots, z_K) \in \mathbb{R}^K$$

Word embedding is any of a set of language modeling and feature learning techniques in natural language processing where words from the vocabulary are mapped to vectors of real numbers.



Softmax:

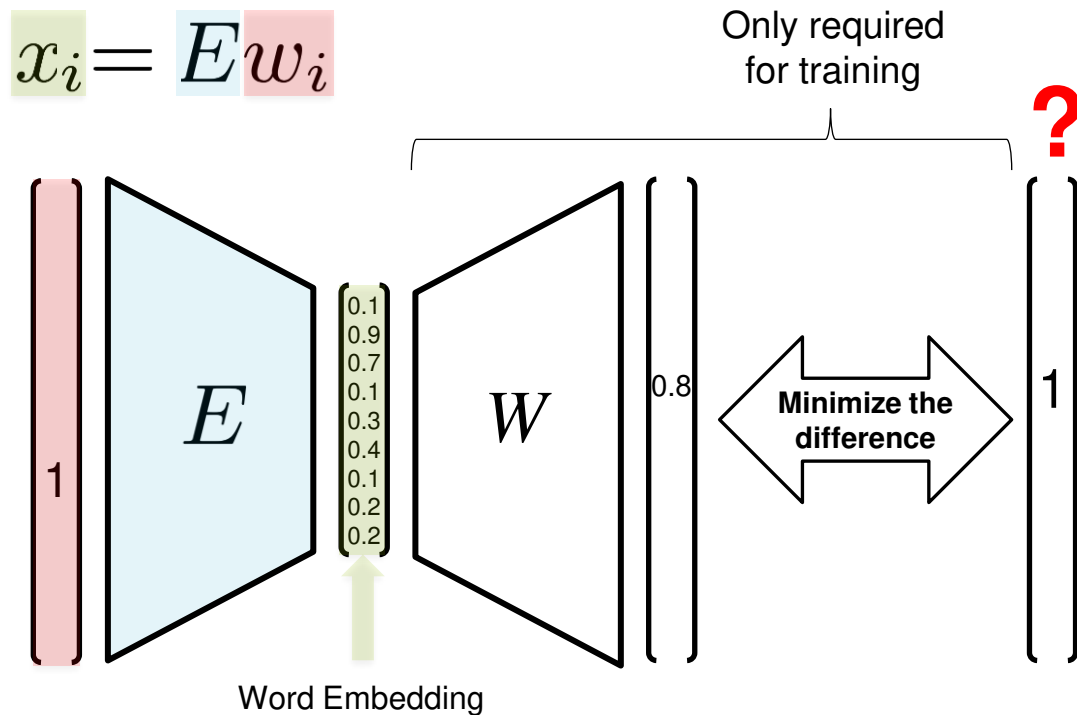
$$\sigma(\mathbf{z})_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}} \text{ for } i = 1, \dots, K \text{ and } \mathbf{z} = (z_1, \dots, z_K) \in \mathbb{R}^K$$

# Word Embedding

## Transform a Word Vector



Word embedding is any of a set of language modeling and feature learning techniques in natural language processing where words from the vocabulary are mapped to vectors of real numbers.

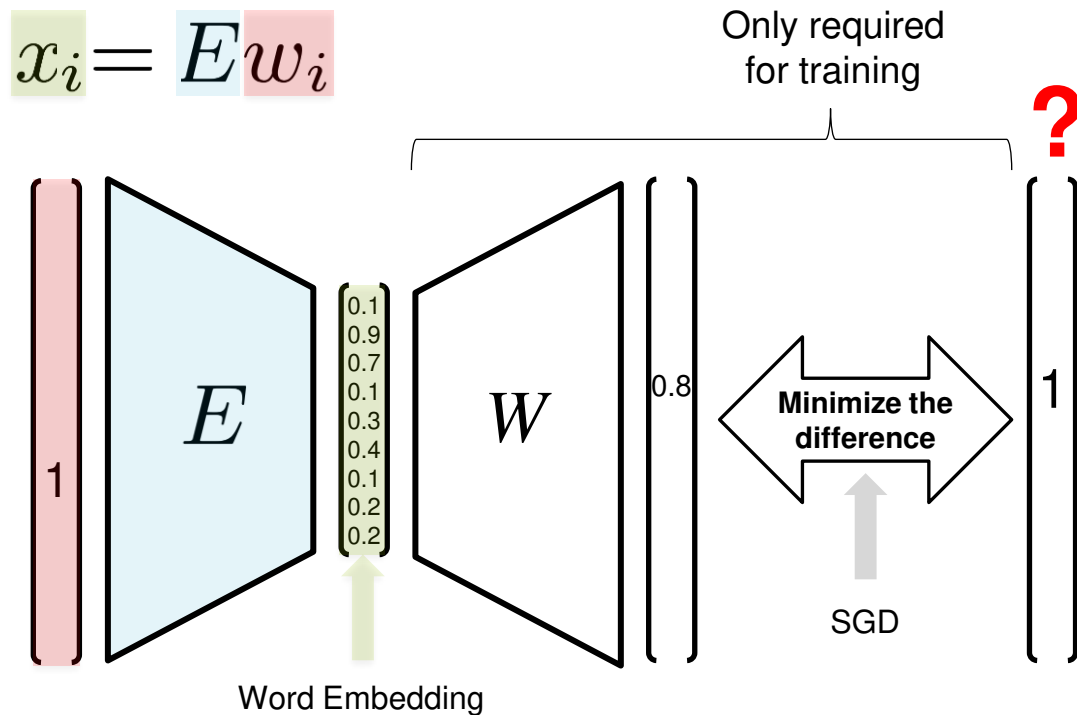


# Word Embedding

## Transform a Word Vector



Word embedding is any of a set of language modeling and feature learning techniques in natural language processing where words from the vocabulary are mapped to vectors of real numbers.

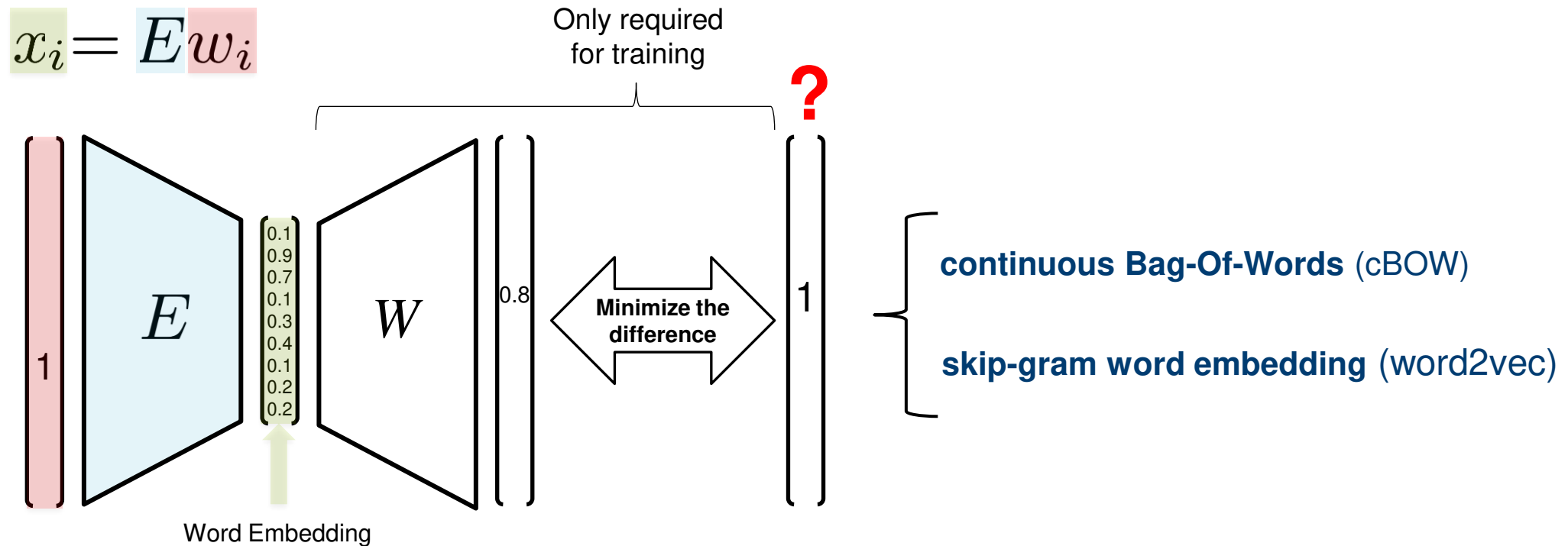


# Word Embedding

## Transform a Word Vector



Word embedding is any of a set of language modeling and feature learning techniques in natural language processing where words from the vocabulary are mapped to vectors of real numbers.



# *Sparse vs. Dense vectors*

## *Differences*



### **tf-idf (or PMI) vectors are**

- Long (length of vocabulary, 20.000 to 50.000)
- Sparse (most elements are zero)

# *Sparse vs. Dense vectors*

## *Differences*

### **tf-idf (or PMI) vectors are**

- Long (length of vocabulary, 20.000 to 50.000)
- Sparse (most elements are zero)

### **Alternative: we want vectors that are**

- Short (length 50-1000)
- Dense (most elements non-zero)



# *Sparse vs. Dense vectors*

## *Question*

### **tf-idf (or PMI) vectors are**

- Long (length of vocabulary, 20.000 to 50.000)
- Sparse (most elements are zero)

### **Alternative: we want vectors that are**

- Short (length 50-1000)
- Dense (most elements non-zero)

**Why (short) dense vectors?**

# Sparse vs. Dense vectors

## Motivation for dense vectors

### tf-idf (or PMI) vectors are

- Long (length of vocabulary, 20.000 to 50.000)
- Sparse (most elements are zero)

### Alternative: we want vectors that are

- Short (length 50-1000)
- Dense (most elements non-zero)

### Why (short) dense vectors?

They may be **easier to use as features** in machine learning (fewer weights to tune)

- Dense vectors may **generalize better** than explicit counts
- Dense vectors may do **better at capturing synonymy**

**They work better in practice**

# Motivation

*Word2vec (Mikolov et al)*

<https://code.google.com/archive/p/word2vec/>

- Popular embedding method
- Very fast to train
- Code available on the web
- Idea: predict rather than count
- Word2vec provides various options.

## **SELF-SUPERVISION**

- No need for human labels
- Bengio et al. (2003); Collobert et al. (2011)

# Word2Vec

*Skip Gram Idea: „Understanding  $\Leftrightarrow$  generating context“*



**Input:** word  
**Output:** its surrounding or *context words*

# Word2Vec

Skip Gram Idea: „Understanding  $\Leftrightarrow$  generating context“



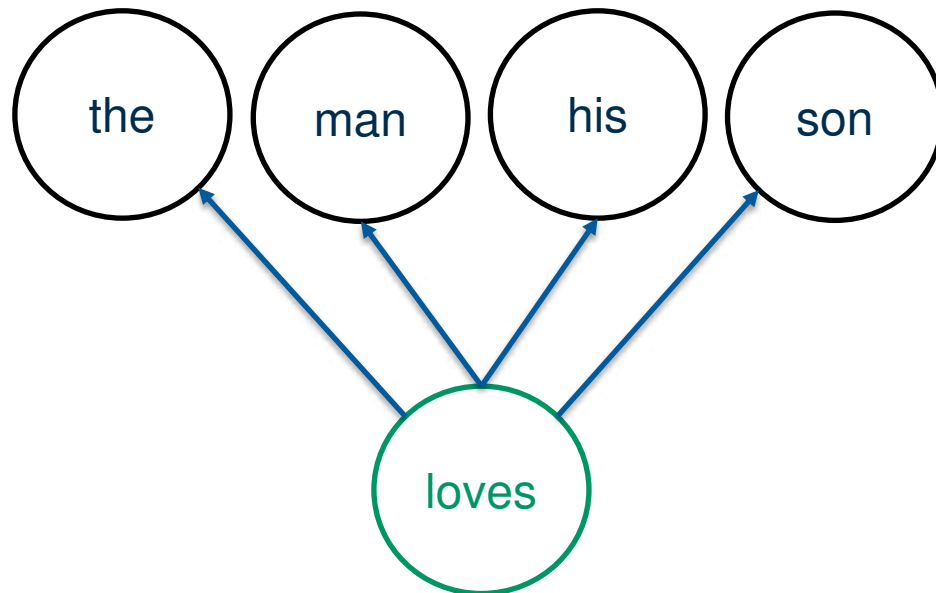
**Input:** word  
**Output:** its surrounding or *context words*

**Example:** “the man loves his son”

Skip Gram Idea: „Understanding  $\Leftrightarrow$  generating context“

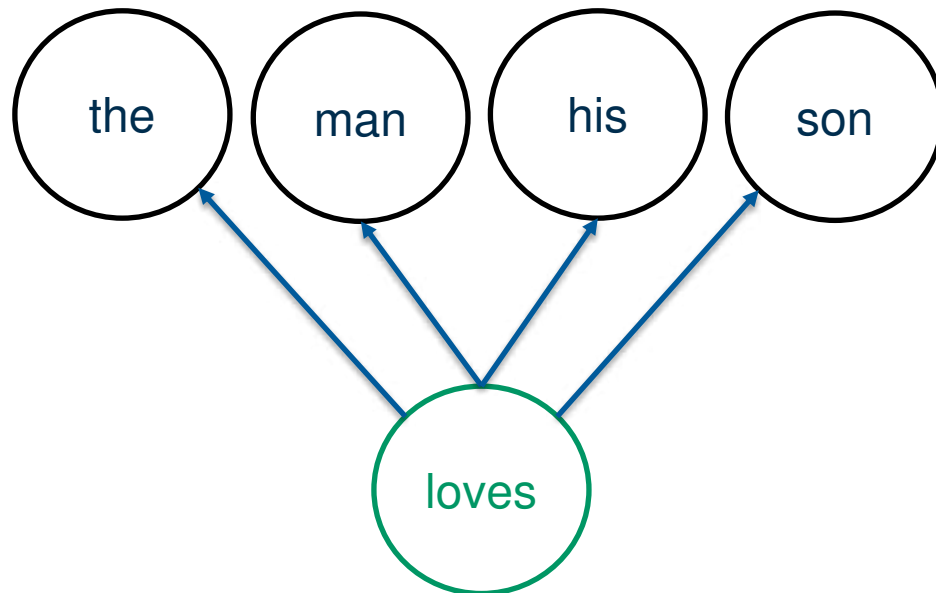
**Input:** word  
**Output:** its surrounding or *context* words

**Example:** “the man loves his son”



**Input:** word  
**Output:** its surrounding or *context* words

**Example:** “the man loves his son”



Skip gram considers conditional probabilities

$$P(\text{"the", "man", "his", "son"} \mid \text{"loves"})$$

# Word2vec: Skip Gram

## Formalization

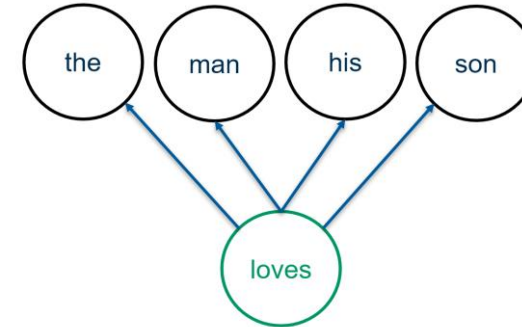
- $T = |X|$  length of text  $X$
- $w^{(t)}$  = word at time step  $t$
- $m$  = context window size
- $\mathbf{v}_i \in \mathbb{R}^d$  = d-dimensional vector, when used as **center word**
- $\mathbf{u}_i \in \mathbb{R}^d$  = d-dimensional vector, when used as **context word**



# Word2vec: Skip Gram

## Formalization

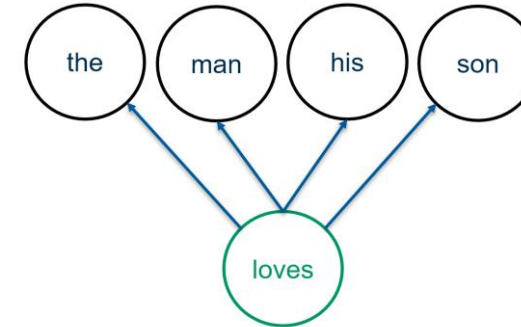
- $T = |X|$  = length of text  $X$
- $w^{(t)}$  = word at time step  $t$
- $m$  = context window size
- $\mathbf{v}_i \in \mathbb{R}^d$  = d-dimensional vector, when used as **center word**
- $\mathbf{u}_i \in \mathbb{R}^d$  = d-dimensional vector, when used as **context word**



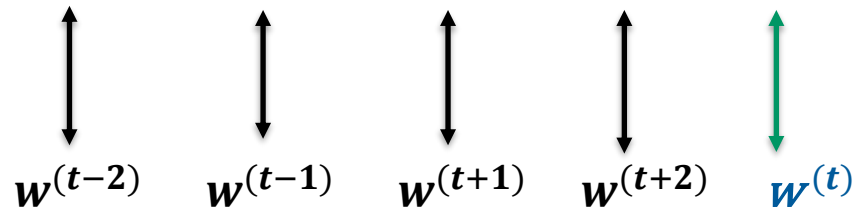
# Word2vec: Skip Gram

## Formalization

- $T = |X|$  length of text  $X$
- $w^{(t)}$  = word at time step  $t$
- $m$  = context window size
- $\mathbf{v}_i \in \mathbb{R}^d$  = d-dimensional vector, when used as **center word**
- $\mathbf{u}_i \in \mathbb{R}^d$  = d-dimensional vector, when used as **context word**



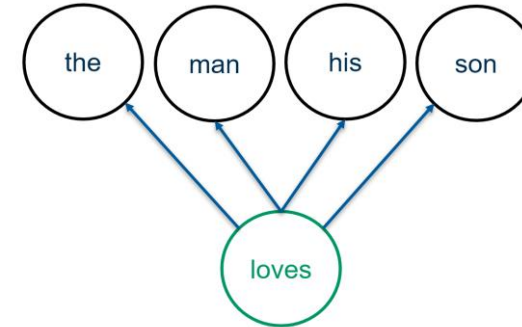
$P(\text{"the", "man", "his", "son" | "loves"})$



# Word2vec: Skip Gram

## Formalization

- $T = |X|$  length of text  $X$
- $w^{(t)}$  = word at time step  $t$
- $m$  = context window size
- $\mathbf{v}_i \in \mathbb{R}^d$  = d-dimensional vector, when used as **center word**
- $\mathbf{u}_i \in \mathbb{R}^d$  = d-dimensional vector, when used as **context word**



$$P(\text{"the", "man", "his", "son" | "loves"}) \stackrel{*}{=} P(\mathbf{w}^{(t-2)} | \mathbf{w}^{(t)}) \cdot P(\mathbf{w}^{(t-1)} | \mathbf{w}^{(t)}) \cdot P(\mathbf{w}^{(t+1)} | \mathbf{w}^{(t)}) \cdot P(\mathbf{w}^{(t+2)} | \mathbf{w}^{(t)})$$

$\mathbf{w}^{(t-2)}$

$\mathbf{w}^{(t-1)}$

$\mathbf{w}^{(t+1)}$

$\mathbf{w}^{(t+2)}$

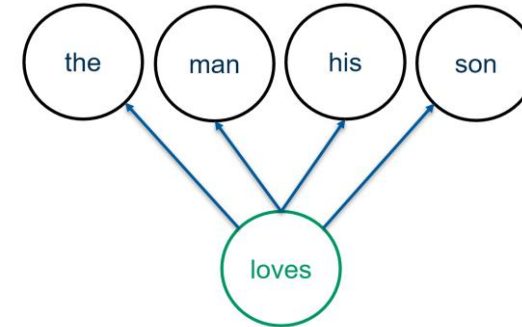
$\mathbf{w}^{(t)}$

*\* context words independently generated given any center word*

# Word2vec: Skip Gram

## Formalization

- $T = |X|$  length of text  $X$
- $w^{(t)}$  = word at time step  $t$
- $m$  = context window size
- $\mathbf{v}_i \in \mathbb{R}^d$  = d-dimensional vector, when used as **center word**
- $\mathbf{u}_i \in \mathbb{R}^d$  = d-dimensional vector, when used as **context word**



$$\begin{aligned}
 P(\text{"the", "man", "his", "son" | "loves"}) &=^* P(w^{(t-2)} | w^{(t)}) \cdot P(w^{(t-1)} | w^{(t)}) \cdot P(w^{(t+1)} | w^{(t)}) \cdot P(w^{(t+2)} | w^{(t)}) \\
 &= \prod_{j=-2, j \neq 0}^2 P(w^{(t+j)} | w^{(t)})
 \end{aligned}$$

$\begin{matrix} \updownarrow & \updownarrow & \updownarrow & \updownarrow & \updownarrow \\ w^{(t-2)} & w^{(t-1)} & w^{(t+1)} & w^{(t+2)} & w^{(t)} \end{matrix}$

*\* context words independently generated given any center word*

\* context words independently generated given any center word

# Word2vec: Skip Gram

## Formalization

- $T = |X|$  length of text  $X$
- $w^{(t)}$  = word at time step  $t$
- $m$  = context window size
- $\mathbf{v}_i \in \mathbb{R}^d$  =  $d$ -dimensional vector, when used as **center word**
- $\mathbf{u}_i \in \mathbb{R}^d$  =  $d$ -dimensional vector, when used as **context word**

## Model parameters

$$\theta = \begin{bmatrix} v_{aardvark} \\ v_a \\ \vdots \\ v_{zebra} \\ u_{aardvark} \\ u_a \\ \vdots \\ u_{zebra} \end{bmatrix} \in \mathbb{R}^{2d|V|}$$

$$P(\text{"the", "man", "his", "son" | "loves"}) = P(w^{(t-2)} | w^{(t)}) \cdot P(w^{(t-1)} | w^{(t)}) \cdot P(w^{(t+1)} | w^{(t)}) \cdot P(w^{(t+2)} | w^{(t)})$$

$$= \prod_{j=-2, j \neq 0}^2 P(w^{(t+j)} | w^{(t)})$$

$\begin{matrix} \updownarrow & \updownarrow & \updownarrow & \updownarrow & \updownarrow \\ w^{(t-2)} & w^{(t-1)} & w^{(t+1)} & w^{(t+2)} & w^{(t)} \end{matrix}$

\* context words independently generated given any center word

# Skip Gram

## Deriving the loss function



X = „This is an example text to learn word2vec algorithm. We can consider sentences like **The man loves his son so much that ...**“  
 $w^{(t)}$

$$\prod_{j=-2, j \neq 0}^2 P(w^{(t+j)} | w^{(t)})$$

# Skip Gram

## Deriving the loss function



X = „This is an example text to learn word2vec algorithm. We can consider sentences like *The **man loves his son so much that ...***“

$w^{(t+1)}$

$$\prod_{j=-2, j \neq 0}^2 P(w^{(t+j)} | w^{(t+1)})$$



# Skip Gram

## Deriving the loss function



X = „This is an example text to learn word2vec algorithm. We can consider sentences like *The man loves his son so much that ...*“  
 $w^{(t+2)}$

$$\prod_{j=-2, j \neq 0}^2 P(w^{(t+j)} | w^{(t+2)})$$

# Skip Gram

## Deriving the loss function



X = „This is an example text to learn word2vec algorithm. We can consider sentences like **The man loves his son so much that ...**“  
 $w^{(t)}$

**Idea:** We want to have good word vectors for the whole text corpus X

# Skip Gram

## Deriving the loss function

$X$  = „This is an example text to learn word2vec algorithm. We can consider sentences like **The man loves his son so much that ...**“  
 $w^{(t)}$

**Idea:** We want to have good word vectors for the whole text corpus  $X$

$$\text{maximize} \quad \prod_{\text{center}} \prod_{\text{context}} P(\text{context} \mid \text{center})$$

# Skip Gram

## Deriving the loss function



X = „This is an example text to learn word2vec algorithm. We can consider sentences like **The man loves his son so much that ...**“  
 $w^{(t)}$

**Idea:** We want to have good word vectors for the whole text corpus X

$$\begin{array}{ccc} \text{maximize} & \prod_{\text{center}} \prod_{\text{context}} P(\text{context} | \text{center}) & \\ \downarrow & & \\ \text{maximize} & \prod_{t=1}^T \prod_{j=-m, j \neq 0}^m P(w^{(t+j)} | w^{(t)}) & \end{array}$$

# Skip Gram

## Deriving the loss function



X = „This is an example text to learn word2vec algorithm. We can consider sentences like **The man loves his son so much that ...**“  
 $w^{(t)}$

**Idea:** We want to have good word vectors for the whole text corpus X

$$\begin{array}{l} \text{maximize} \\ \prod_{t=1}^T \prod_{j=-m, j \neq 0}^m P(w^{(t+j)} | w^{(t)}) \\ \downarrow \\ \text{minimize} \\ - \prod_{t=1}^T \prod_{j=-m, j \neq 0}^m P(w^{(t+j)} | w^{(t)}) \end{array}$$

# Skip Gram

## Deriving the loss function



$X$  = „This is an example text to learn word2vec algorithm. We can consider sentences like **The man loves his son so much that ...**“  
 $w^{(t)}$

**Idea:** We want to have good word vectors for the whole text corpus  $X$

$$\begin{array}{l} \text{minimize} \\ \log \downarrow \\ \text{minimize} \end{array} \quad - \prod_{t=1}^T \prod_{j=-m, j \neq 0}^m P(w^{(t+j)} | w^{(t)})$$
$$- \sum_{t=1}^T \sum_{j=-m, j \neq 0}^m \log P(w^{(t+j)} | w^{(t)})$$

NLLL:= **N**egative **L**og **L**ikelihood **L**oss

# Skip Gram

## Deriving the loss function



X = „This is an example text to learn word2vec algorithm. We can consider sentences like **The man loves his son so much that ...**“  
 $w^{(t)}$

### Questions:

- How to compute the probabilities?
- And how are probabilities and word vectors related?

$$\text{minimize} \quad - \sum_{t=1}^T \sum_{j=-m, j \neq 0}^m \log P(w^{(t+j)} | w^{(t)})$$



Let  $\mathbf{x} \in \mathbb{R}^n$ . The *softmax* function is defined as

$$\text{softmax}(\mathbf{x})_j = \frac{\exp(x_j)}{\sum_{i=0}^{n-1} \exp(x_i)}, \quad j = 0, \dots, n-1$$

## Question:

It is a function

$$\text{softmax}: \mathbb{R}^? \rightarrow (a, b)^?$$





Let  $x \in \mathbb{R}^n$ . The function  $\text{softmax} : \mathbb{R}^n \rightarrow (0,1)^n$  is defined as

$$\text{softmax}(x)_j = \frac{\exp(x_j)}{\sum_{i=0}^{n-1} \exp(x_i)}, \quad j = 0, \dots, n-1$$

# Softmax Function

And probabilities



Let  $x \in \mathbb{R}^n$ . The function  $\text{softmax} : \mathbb{R}^n \rightarrow (0,1)^n$  is defined as

$$\text{softmax}(x)_j = \frac{\exp(x_j)}{\sum_{i=0}^{n-1} \exp(x_i)}, \quad j = 0, \dots, n-1$$

**Question:** How to compute the probabilities? And what's their relation to word vectors?

# Softmax Function

And probabilities

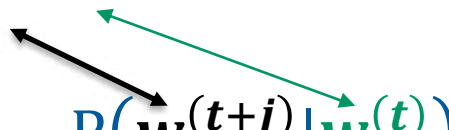
Let  $\mathbf{x} \in \mathbb{R}^n$ . The function  $\text{softmax} : \mathbb{R}^n \rightarrow (0,1)^n$  is defined as

$$\text{softmax}(\mathbf{x})_j = \frac{\exp(x_j)}{\sum_{i=0}^{n-1} \exp(x_i)}, \quad j = 0, \dots, n-1$$

**Question:** How to compute the probabilities? And what's their relation to word vectors?

Set  $x_j = \mathbf{u}_o^T \mathbf{v}_c$  and assume  $n = |V| \Leftrightarrow$  vocabulary size

$\Rightarrow$

$$P(\mathbf{w}^{(t+j)} | \mathbf{w}^{(t)}) = \text{softmax}(\mathbf{u}_o^T \mathbf{v}_c) = \frac{\exp(\mathbf{u}_o^T \mathbf{v}_c)}{\sum_{i=0}^{n-1} \exp(\mathbf{u}_i^T \mathbf{v}_c)}$$


# Softmax Function

And probabilities

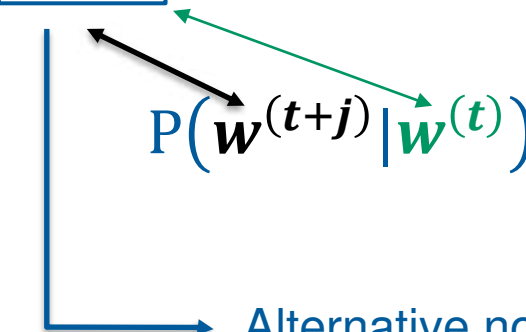
Let  $\mathbf{x} \in \mathbb{R}^n$ . The function  $\text{softmax} : \mathbb{R}^n \rightarrow (0,1)^n$  is defined as

$$\text{softmax}(\mathbf{x})_j = \frac{\exp(x_j)}{\sum_{i=0}^{n-1} \exp(x_i)}, \quad j = 0, \dots, n-1$$

**Question:** How to compute the probabilities? And what's their relation to word vectors?

Set  $x_j = \boxed{\mathbf{u}_o^T \mathbf{v}_c}$  and assume  $n = |V| \Leftrightarrow$  vocabulary size

=>

$$P(\mathbf{w}^{(t+j)} | \mathbf{w}^{(t)}) = \text{softmax}(\mathbf{u}_o^T \mathbf{v}_c) = \frac{\exp(\mathbf{u}_o^T \mathbf{v}_c)}{\sum_{i=0}^{n-1} \exp(\mathbf{u}_i^T \mathbf{v}_c)}$$


Alternative notations for dot product:  $\mathbf{u}_o^T \mathbf{v}_c = \langle \mathbf{u}_o, \mathbf{v}_c \rangle = \mathbf{u}_o \cdot \mathbf{v}_c = \sum_{k=1}^d \mathbf{u}_{o_k} \mathbf{v}_{c_k}$

# Softmax Function

And probabilities

Let  $\mathbf{x} \in \mathbb{R}^n$ . The function  $\text{softmax} : \mathbb{R}^n \rightarrow (0,1)^n$  is defined as

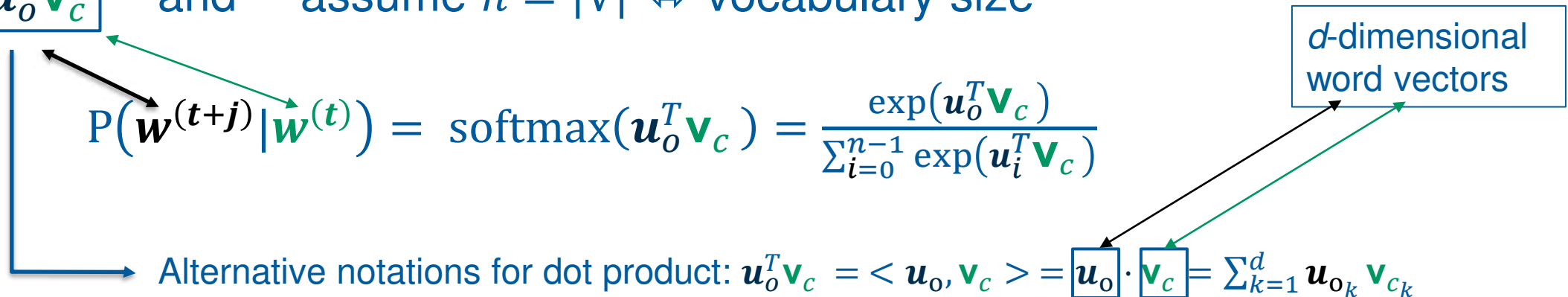
$$\text{softmax}(\mathbf{x})_j = \frac{\exp(x_j)}{\sum_{i=0}^{n-1} \exp(x_i)}, \quad j = 0, \dots, n-1$$

**Question:** How to compute the probabilities? And what's their relation to word vectors?

Set  $x_j = \mathbf{u}_o^T \mathbf{v}_c$  and assume  $n = |V| \Leftrightarrow$  vocabulary size

$\Rightarrow$

$$P(\mathbf{w}^{(t+j)} | \mathbf{w}^{(t)}) = \text{softmax}(\mathbf{u}_o^T \mathbf{v}_c) = \frac{\exp(\mathbf{u}_o^T \mathbf{v}_c)}{\sum_{i=0}^{n-1} \exp(\mathbf{u}_i^T \mathbf{v}_c)}$$



# Softmax Function

And probabilities

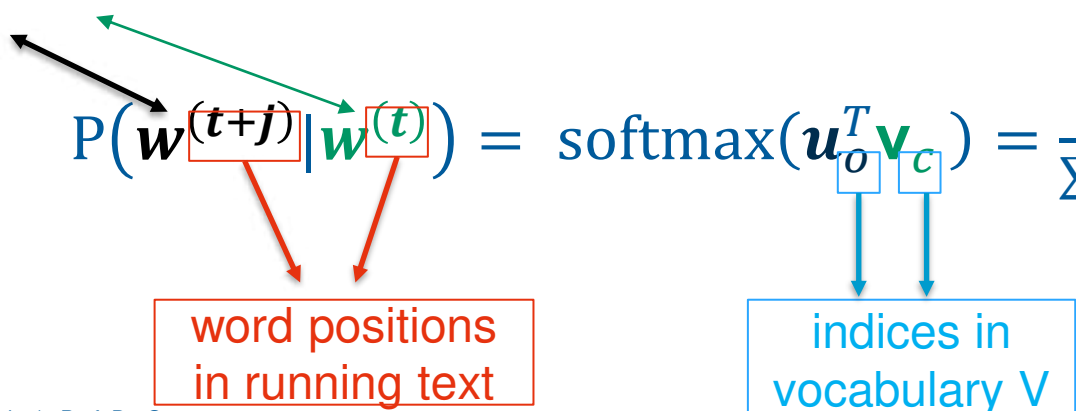
Let  $\mathbf{x} \in \mathbb{R}^n$ . The function  $\text{softmax} : \mathbb{R}^n \rightarrow (0,1)^n$  is defined as

$$\text{softmax}(\mathbf{x})_j = \frac{\exp(x_j)}{\sum_{i=0}^{n-1} \exp(x_i)}, \quad j = 0, \dots, n-1$$

**Question:** How to compute the probabilities? And what's their relation to word vectors?

Set  $x_j = \mathbf{u}_o^T \mathbf{v}_c$  and assume  $n = |V| \Leftrightarrow$  vocabulary size

=>

$$P(\mathbf{w}^{(t+j)} | \mathbf{w}^{(t)}) = \text{softmax}(\mathbf{u}_o^T \mathbf{v}_c) = \frac{\exp(\mathbf{u}_o^T \mathbf{v}_c)}{\sum_{i=0}^{n-1} \exp(\mathbf{u}_i^T \mathbf{v}_c)}$$


word positions in running text

indices in vocabulary V

# Skip Gram

## Deriving the loss function

$X =$  „This is an example text to learn word2vec algorithm. We can consider sentences like **The man loves his son so much that ...**“  
 $w^{(t)}$

**Question:** How to compute the probabilities?

$$\text{minimize} \quad - \sum_{t=1}^T \sum_{j=-m, j \neq 0}^m \log P(w^{(t+j)} | w^{(t)})$$



$\mathbf{v}_c :=$  center word with index  $c \in V$   
 $\mathbf{u}_i, \mathbf{u}_o :=$  context words with indices  $i, o \in V$

$$P(w_o | w_c) = \frac{\exp(\mathbf{u}_o^\top \mathbf{v}_c)}{\sum_{i \in V} \exp(\mathbf{u}_i^\top \mathbf{v}_c)},$$

Vocabulary  $V = \{0, 1, \dots, |V| - 1\}$

# Skip Gram: Gradients

Minimizing the loss function

$$\text{minimize} \quad - \sum_{t=1}^T \sum_{j=-m, j \neq 0}^m \log P(\mathbf{w}^{(t+j)} | \mathbf{w}^{(t)})$$

Vocabulary  $V = \{0, 1, \dots, |V| - 1\}$

$\mathbf{v}_c :=$  center word with index  $c \in V$   
 $\mathbf{u}_i :=$  context word with index  $i \in V$



# Skip Gram: Gradients

Minimizing the loss function

minimize

$$L(\theta) = - \sum_{t=1}^T \sum_{j=-m, j \neq 0}^m \log P(\mathbf{w}^{(t+j)} | \mathbf{w}^{(t)})$$

Theta holds all variables (center and context word vectors)



Vocabulary  $V = \{0, 1, \dots, |V| - 1\}$

$\mathbf{v}_c$  := center word with index  $c \in V$   
 $\mathbf{u}_i$  := context word with index  $i \in V$

# Skip Gram: Gradients

Minimizing the loss function

minimize

$$L(\theta) = - \sum_{t=1}^T \sum_{j=-m, j \neq 0}^m \log P(\mathbf{w}^{(t+j)} | \mathbf{w}^{(t)})$$

Theta holds all variables (center and context word vectors)



Vocabulary  $V = \{0, 1, \dots, |V| - 1\}$

$\mathbf{v}_c :=$  center word with index  $c \in V$   
 $\mathbf{u}_i :=$  context word with index  $i \in V$

# Skip Gram: Gradients

Minimizing the loss function



minimize

$$L(\theta) = - \sum_{t=1}^T \sum_{j=-m, j \neq 0}^m \log P(\mathbf{w}^{(t+j)} | \mathbf{w}^{(t)})$$

Theta holds all variables (center and context word vectors)

$$P(w_o | w_c) = \frac{\exp(\mathbf{u}_o^\top \mathbf{v}_c)}{\sum_{i \in \mathcal{V}} \exp(\mathbf{u}_i^\top \mathbf{v}_c)},$$

Question: What is the log of P(...) ?

Vocabulary  $V = \{0, 1, \dots, |V| - 1\}$

$\mathbf{v}_c :=$  center word with index  $c \in V$   
 $\mathbf{u}_i :=$  context word with index  $i \in V$

# Skip Gram: Gradients

Minimizing the loss function

minimize

$$P(w_o | w_c) = \frac{\exp(\mathbf{u}_o^\top \mathbf{v}_c)}{\sum_{i \in \mathcal{V}} \exp(\mathbf{u}_i^\top \mathbf{v}_c)},$$

$$\log P(w_o | w_c) = \mathbf{u}_o^\top \mathbf{v}_c - \log \left( \sum_{i \in \mathcal{V}} \exp(\mathbf{u}_i^\top \mathbf{v}_c) \right)$$

↗ Theta holds all variables (word vectors)

$$L(\theta) = - \sum_{t=1}^T \sum_{j=-m, j \neq 0}^m \log P(\mathbf{w}^{(t+j)} | \mathbf{w}^{(t)})$$

**Question:** What is the log of P(...)?

Vocabulary  $V = \{0, 1, \dots, |V| - 1\}$

$\mathbf{v}_c :=$  center word with index  $c \in V$   
 $\mathbf{u}_i :=$  context word with index  $i \in V$



# Skip Gram: Gradients

Minimizing the loss function

minimize

$$L(\theta) = - \sum_{t=1}^T \sum_{j=-m, j \neq 0}^m \log P(\mathbf{w}^{(t+j)} | \mathbf{w}^{(t)})$$

Theta holds all variables (word vectors)

$$P(w_o | w_c) = \frac{\exp(\mathbf{u}_o^\top \mathbf{v}_c)}{\sum_{i \in \mathcal{V}} \exp(\mathbf{u}_i^\top \mathbf{v}_c)},$$

**Question:** Which derivatives do we need to compute?

Vocabulary  $V = \{0, 1, \dots, |V| - 1\}$

$\mathbf{v}_c :=$  center word with index  $c \in V$   
 $\mathbf{u}_i :=$  context word with index  $i \in V$

# Skip Gram: Gradients

Minimizing the loss function



Theta holds all variables (word vectors)

$L(\theta)$

minimize

$$-\sum_{t=1}^T \sum_{j=-m, j \neq 0}^m \log P(\mathbf{w}^{(t+j)} | \mathbf{w}^{(t)})$$

$$P(w_o | w_c) = \frac{\exp(\mathbf{u}_o^\top \mathbf{v}_c)}{\sum_{i \in \mathcal{V}} \exp(\mathbf{u}_i^\top \mathbf{v}_c)},$$

Compute  $\nabla_{\theta} L(\theta)$ , where  $\theta = (u_0, \dots, u_{|V|-1}, \mathbf{v}_0, \dots, \mathbf{v}_{|V|-1})$ ,  $u_i, \mathbf{v}_j \in \mathbb{R}^d$

Vocabulary  $V = \{0, 1, \dots, |V| - 1\}$

$\mathbf{v}_c :=$  center word with index  $c \in V$   
 $\mathbf{u}_i :=$  context word with index  $i \in V$

# Skip Gram: Gradients

Minimizing the loss function

minimize

$$L(\theta) \quad \text{Theta holds all variables (word vectors)}$$

$$- \sum_{t=1}^T \sum_{j=-m, j \neq 0}^m \log P(\mathbf{w}^{(t+j)} | \mathbf{w}^{(t)})$$

$$P(w_o | w_c) = \frac{\exp(\mathbf{u}_o^\top \mathbf{v}_c)}{\sum_{i \in \mathcal{V}} \exp(\mathbf{u}_i^\top \mathbf{v}_c)},$$

Compute  $\nabla_{\theta} L(\theta)$ , where  $\theta = (u_0, \dots, u_{|V|-1}, \mathbf{v}_0, \dots, \mathbf{v}_{|V|-1})$ ,  $u_i, \mathbf{v}_j \in \mathbb{R}^d$

Derivative w.r.t. center word vector parameters:

$$\begin{aligned} \frac{\partial \log P(w_o | w_c)}{\partial \mathbf{v}_c} &= \mathbf{u}_o - \frac{\sum_{j \in \mathcal{V}} \exp(\mathbf{u}_j^\top \mathbf{v}_c) \mathbf{u}_j}{\sum_{i \in \mathcal{V}} \exp(\mathbf{u}_i^\top \mathbf{v}_c)} \\ &= \mathbf{u}_o - \sum_{j \in \mathcal{V}} \left( \frac{\exp(\mathbf{u}_j^\top \mathbf{v}_c)}{\sum_{i \in \mathcal{V}} \exp(\mathbf{u}_i^\top \mathbf{v}_c)} \right) \mathbf{u}_j \\ &= \mathbf{u}_o - \sum_{j \in \mathcal{V}} P(w_j | w_c) \mathbf{u}_j. \end{aligned}$$

Vocabulary  $V = \{0, 1, \dots, |V| - 1\}$

$\mathbf{v}_c :=$  center word with index  $c \in V$   
 $\mathbf{u}_i :=$  context word with index  $i \in V$

# Skip Gram: Gradients

Minimizing the loss function

minimize

$$L(\theta) \quad \text{Theta holds all variables (word vectors)}$$

$$- \sum_{t=1}^T \sum_{j=-m, j \neq 0}^m \log P(\mathbf{w}^{(t+j)} | \mathbf{w}^{(t)})$$

$$P(w_o | w_c) = \frac{\exp(\mathbf{u}_o^\top \mathbf{v}_c)}{\sum_{i \in \mathcal{V}} \exp(\mathbf{u}_i^\top \mathbf{v}_c)},$$

Compute  $\nabla_{\theta} L(\theta)$ , where  $\theta = (u_0, \dots, u_{|V|-1}, \mathbf{v}_0, \dots, \mathbf{v}_{|V|-1})$ ,  $u_i, \mathbf{v}_j \in \mathbb{R}^d$

Derivative w.r.t. center word vector parameters:

$$\begin{aligned} \frac{\partial \log P(w_o | w_c)}{\partial \mathbf{v}_c} &= \mathbf{u}_o - \frac{\sum_{j \in \mathcal{V}} \exp(\mathbf{u}_j^\top \mathbf{v}_c) \mathbf{u}_j}{\sum_{i \in \mathcal{V}} \exp(\mathbf{u}_i^\top \mathbf{v}_c)} \\ &= \mathbf{u}_o - \sum_{j \in \mathcal{V}} \left( \frac{\exp(\mathbf{u}_j^\top \mathbf{v}_c)}{\sum_{i \in \mathcal{V}} \exp(\mathbf{u}_i^\top \mathbf{v}_c)} \right) \mathbf{u}_j \\ &= \mathbf{u}_o - \underbrace{\sum_{j \in \mathcal{V}} P(w_j | w_c) \mathbf{u}_j}_{\text{Average over all context vectors, weighted by their respective probability}} \end{aligned}$$

Average over all context vectors, weighted by their respective probability

Vocabulary  $V = \{0, 1, \dots, |V| - 1\}$

$\mathbf{v}_c :=$  center word with index  $c \in V$   
 $\mathbf{u}_i :=$  context word with index  $i \in V$



# Skip Gram: Gradients

Minimizing the loss function

Theta holds all variables (word vectors)

minimize

$$L(\theta) = - \sum_{t=1}^T \sum_{j=-m, j \neq 0}^m \log P(\mathbf{w}^{(t+j)} | \mathbf{w}^{(t)})$$

$$P(w_o | w_c) = \frac{\exp(\mathbf{u}_o^\top \mathbf{v}_c)}{\sum_{i \in \mathcal{V}} \exp(\mathbf{u}_i^\top \mathbf{v}_c)},$$

Compute  $\nabla_{\theta} L(\theta)$ , where  $\theta = (u_0, \dots, u_{|V|-1}, \mathbf{v}_0, \dots, \mathbf{v}_{|V|-1})$ ,  $u_i, \mathbf{v}_j \in \mathbb{R}^d$

Derivative w.r.t. center word vector parameters:

$$\begin{aligned} \frac{\partial \log P(w_o | w_c)}{\partial \mathbf{v}_c} &= \mathbf{u}_o - \frac{\sum_{j \in \mathcal{V}} \exp(\mathbf{u}_j^\top \mathbf{v}_c) \mathbf{u}_j}{\sum_{i \in \mathcal{V}} \exp(\mathbf{u}_i^\top \mathbf{v}_c)} \\ &= \mathbf{u}_o - \sum_{j \in \mathcal{V}} \left( \frac{\exp(\mathbf{u}_j^\top \mathbf{v}_c)}{\sum_{i \in \mathcal{V}} \exp(\mathbf{u}_i^\top \mathbf{v}_c)} \right) \mathbf{u}_j \\ &= \mathbf{u}_o - \underbrace{\sum_{j \in \mathcal{V}} P(w_j | w_c) \mathbf{u}_j}_{\text{Average over all context vectors, weighted by their respective probability}} \end{aligned}$$

**Exercise:** derivative w.r.t. *context word vectors*?

Average over all context vectors, weighted by their respective probability

Vocabulary  $V = \{0, 1, \dots, |V| - 1\}$

$\mathbf{v}_c :=$  center word with index  $c \in V$   
 $\mathbf{u}_i :=$  context word with index  $i \in V$

# Gradient Descent

## Idea



Update equation

$$\theta^{new} = \theta^{old} - \alpha \nabla_{\theta^{old}} L(\theta)$$

# Gradient Descent

## Idea



Update equation

$$\theta^{new} = \theta^{old} - \alpha \underbrace{\nabla_{\theta^{old}} L(\theta)}_{\text{Gradient of loss function}}$$

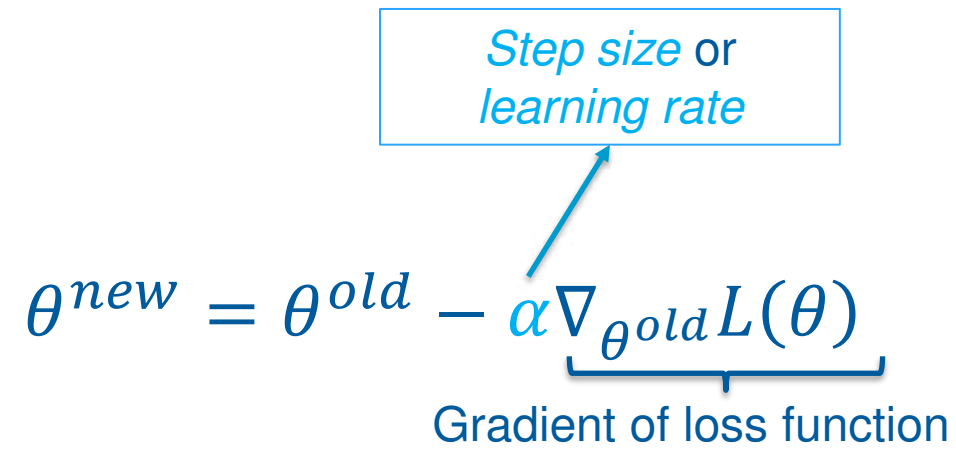
# Gradient Descent

Idea

Update equation

$$\theta^{new} = \theta^{old} - \alpha \underbrace{\nabla_{\theta^{old}} L(\theta)}_{\text{Gradient of loss function}}$$

Step size or learning rate



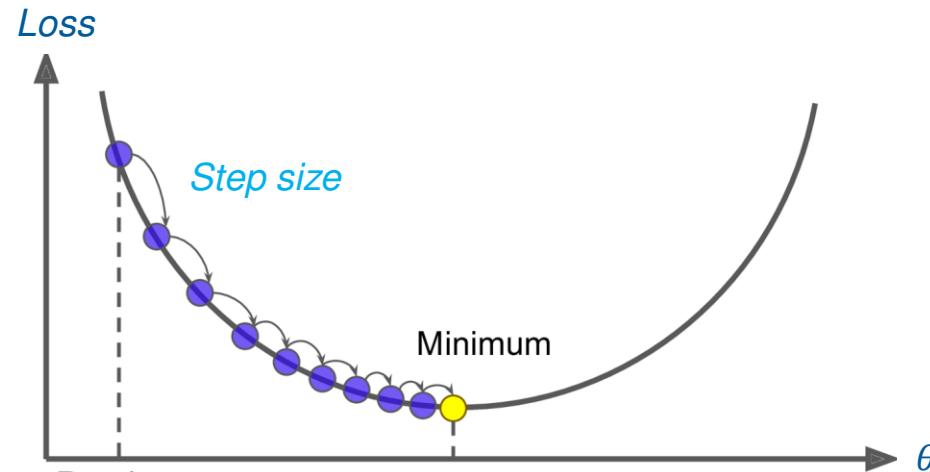
# Gradient Descent

Idea

Update equation

$$\theta^{new} = \theta^{old} - \underbrace{\alpha \nabla_{\theta^{old}} L(\theta)}_{\text{Gradient of loss function}}$$

Step size or  
learning rate



$$\theta = \begin{bmatrix} v_{aardvark} \\ v_a \\ \vdots \\ v_{zebra} \\ u_{aardvark} \\ u_a \\ \vdots \\ u_{zebra} \end{bmatrix}$$

# Word2Vec: Skip Gram Summary

## Idea



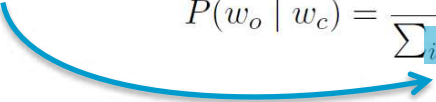
**Problem:**  $L(\theta)$  is a function where softmax is computed over all words (and corresponding windows) in vocabulary

Further Reading: [https://d2l.ai/chapter\\_natural-language-processing-pretraining/approx-training.html](https://d2l.ai/chapter_natural-language-processing-pretraining/approx-training.html)

# Word2Vec: Skip Gram Summary

## Idea

**Problem:**  $L(\theta)$  is a function where softmax is computed over all words (and corresponding windows) in vocabulary

$$P(w_o | w_c) = \frac{\exp(\mathbf{u}_o^\top \mathbf{v}_c)}{\sum_{i \in \mathcal{V}} \exp(\mathbf{u}_i^\top \mathbf{v}_c)},$$


Further Reading: [https://d2l.ai/chapter\\_natural-language-processing-pretraining/approx-training.html](https://d2l.ai/chapter_natural-language-processing-pretraining/approx-training.html)

# Word2Vec: Skip Gram Summary

## Idea

**Problem:**  $L(\theta)$  is a function where softmax is computed over **all words** (and corresponding windows) **in vocabulary**

$$P(w_o | w_c) = \frac{\exp(\mathbf{u}_o^\top \mathbf{v}_c)}{\sum_{i \in \mathcal{V}} \exp(\mathbf{u}_i^\top \mathbf{v}_c)},$$

**Solution:**

- Negative Sampling
- Hierarchical Sampling

Computationally expensive if vocabulary  $V$  is large

... in practice:  $V$  **is** very large!

Further Reading: [https://d2l.ai/chapter\\_natural-language-processing-pretraining/approx-training.html](https://d2l.ai/chapter_natural-language-processing-pretraining/approx-training.html)



### Objective Function

Ein Elektroauto ist ein Automobil mit elektrischem Antrieb.

Objective Function :=  $\max \prod_{\text{center}} \prod_{\text{context}} P(\text{context} | \text{center})$

Softmax:

$$\sigma(\mathbf{z})_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}} \text{ for } i = 1, \dots, K \text{ and } \mathbf{z} = (z_1, \dots, z_K) \in \mathbb{R}^K$$

### Objective Function

Ein Elektroauto ist ein **Automobil** mit elektrischem Antrieb.

$$\text{Objective Function} := \max \prod_{\text{center}} \prod_{\text{context}} P(\text{context} | \text{center})$$

$$\min \left( -\frac{1}{N} \sum_{\text{center}} \sum_{\text{context}} \log(P(\text{context} | \text{center})) \right) := \text{Commonly used Objective Function}$$

Softmax:

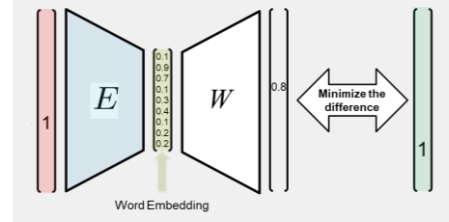
$$\sigma(\mathbf{z})_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}} \text{ for } i = 1, \dots, K \text{ and } \mathbf{z} = (z_1, \dots, z_K) \in \mathbb{R}^K$$

### Objective Function

Ein Elektroauto ist ein **Automobil** mit elektrischem Antrieb.

$$\text{Objective Function} := \max \prod_{\text{center}} \prod_{\text{context}} P(\text{context} | \text{center})$$

$$\min \left( -\frac{1}{N} \sum_{\text{center}} \sum_{\text{context}} \log(P(\text{context} | \text{center})) \right) := \text{Commonly used Objective Function}$$



Softmax:

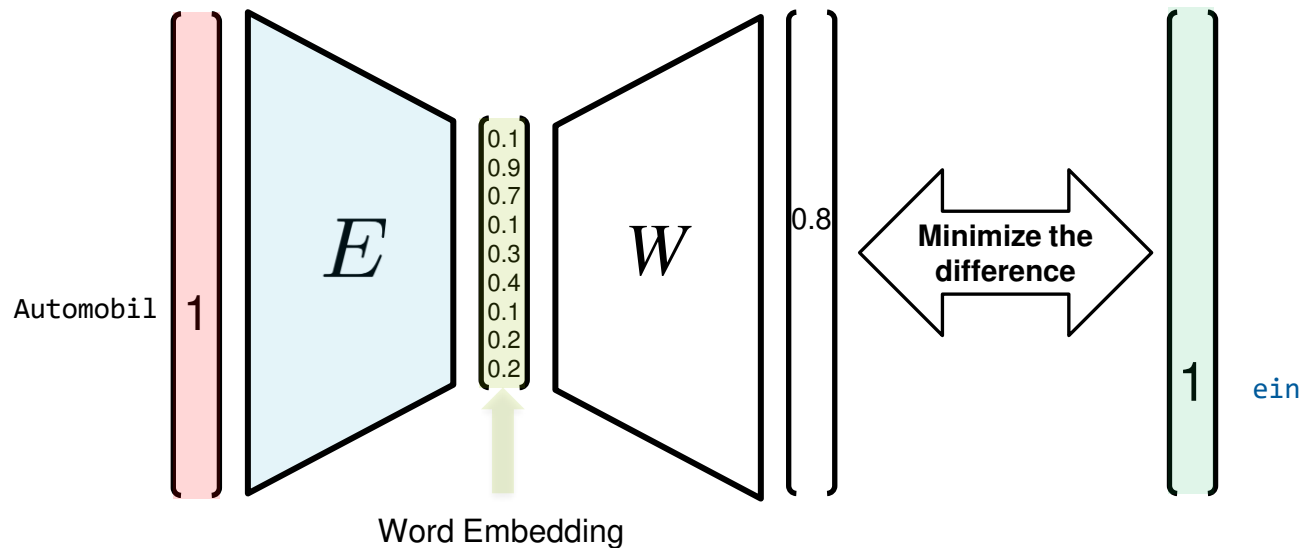
$$\sigma(\mathbf{z})_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}} \text{ for } i = 1, \dots, K \text{ and } \mathbf{z} = (z_1, \dots, z_K) \in \mathbb{R}^K$$

# Word Embedding

## Skip-Gram Word Embedding

Training is computational expensive.

Ein Elektroauto ist ein Automobil mit elektrischem Antrieb.



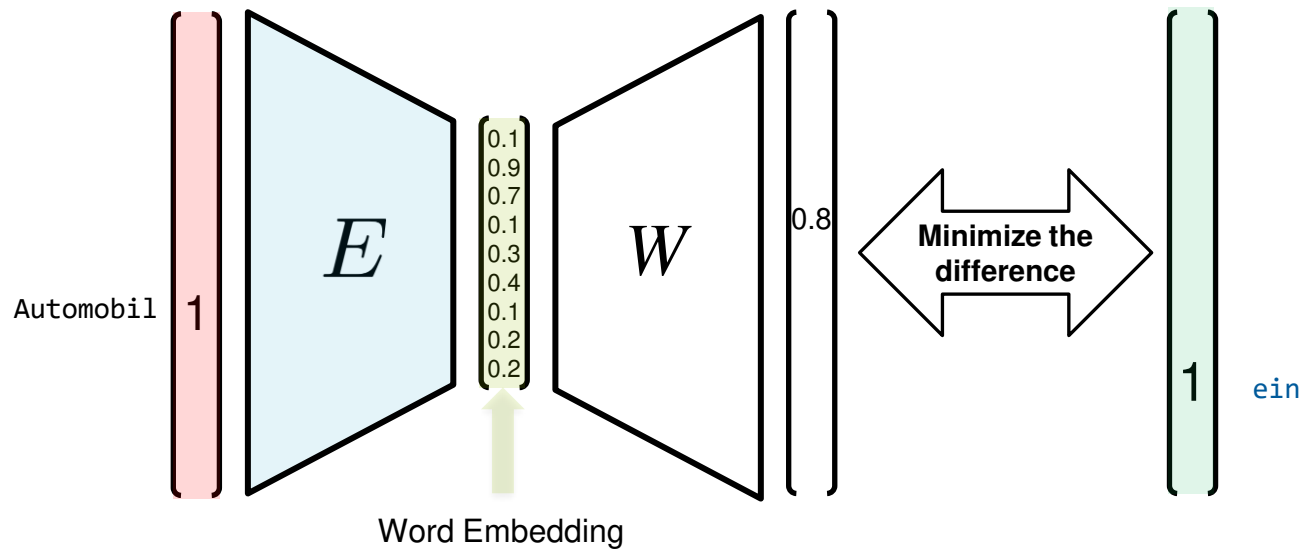
- Any pair of target/context word

# Word Embedding

## Skip-Gram Word Embedding

Training is computational expensive.

Ein Elektroauto ist ein Automobil mit elektrischem Antrieb.



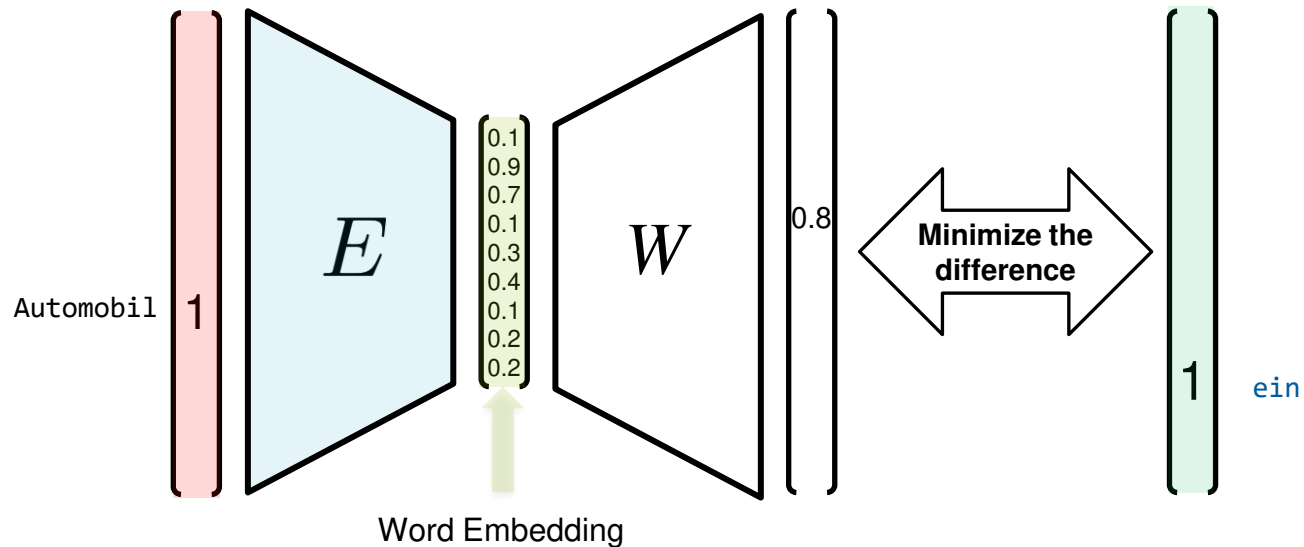
- Any pair of target/context word
- Solutions:
- Sampling

# Word Embedding

## Skip-Gram Word Embedding

Training is computational expensive.

Ein Elektroauto ist ein Automobil mit elektrischem Antrieb.



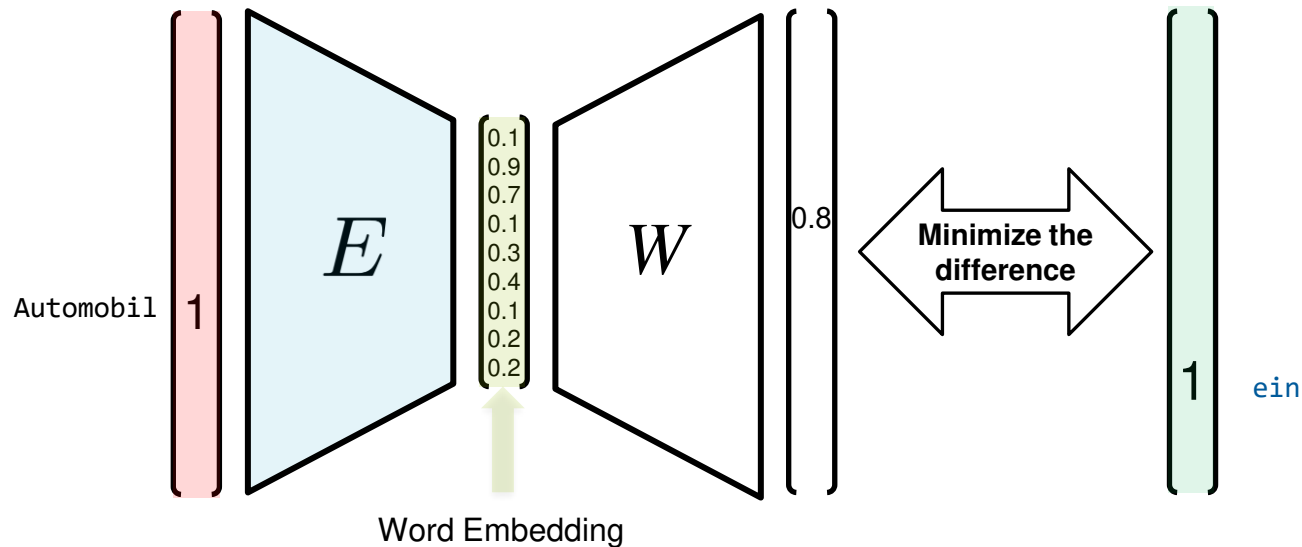
- Any pair of target/context word
- Solutions:
  - Sampling
  - Softmax over vocabulary

# Word Embedding

## Skip-Gram Word Embedding

Training is computational expensive.

Ein Elektroauto ist ein Automobil mit elektrischem Antrieb.



■ Any pair of target/context word

Solutions:

■ Sampling

■ Softmax over vocabulary

Solution:

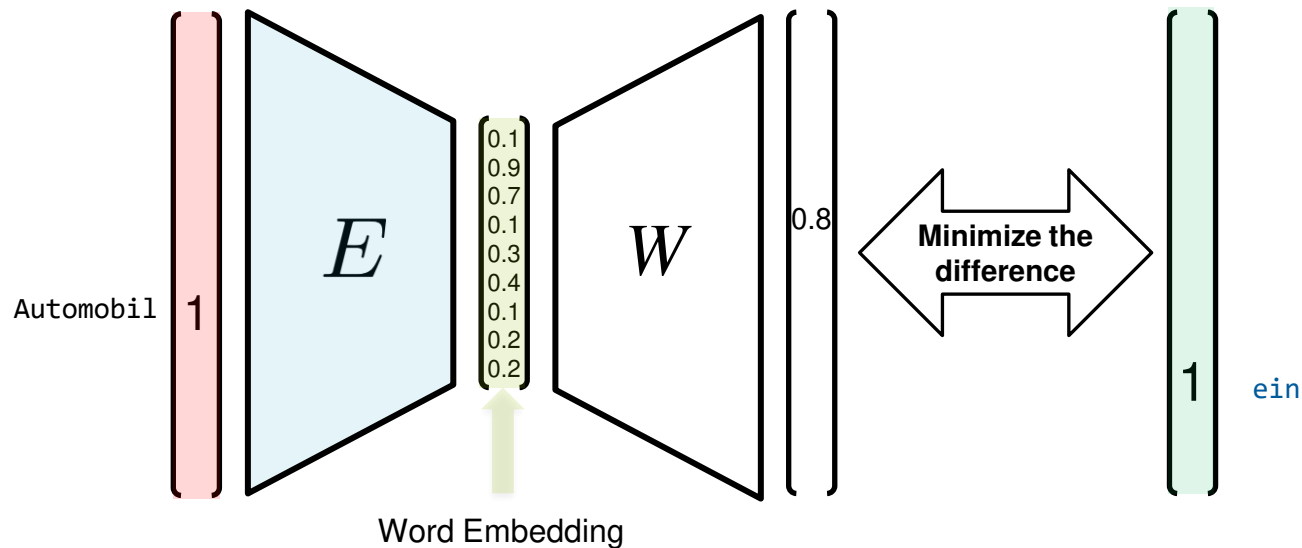
■ Hierarchical Softmax

# Word Embedding

## Skip-Gram Word Embedding

Training is computational expensive.

Ein Elektroauto ist ein Automobil mit elektrischem Antrieb.



- Any pair of target/context word

Solutions:

- Sampling

- Softmax over vocabulary

Solution:

- Hierarchical Softmax
- Noise Contrastive Estimation
- Negative Sampling

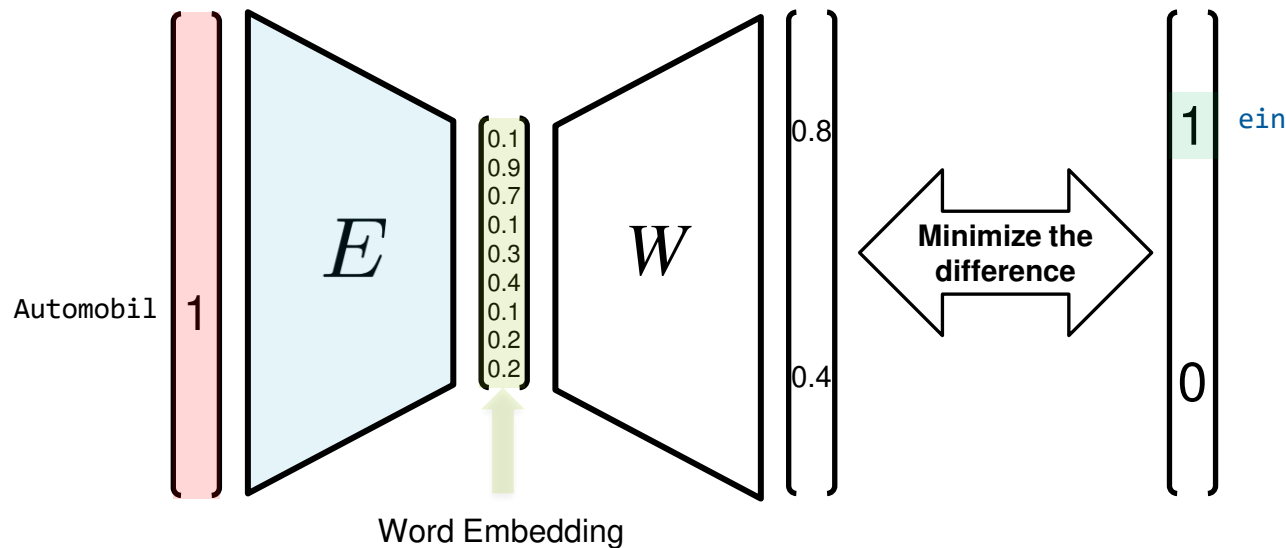


# Word Embedding

## Training Skip-Gram Word Embedding with Negative Sampling

Training is computational expensive.

Ein Elektroauto ist **ein** **Automobil** mit elektrischem Antrieb.



### Move to binary classification:

- Replace Softmax by Sigmoid
- Train with positive and negative samples

Context	Target		
ein	Automobil	:= 1	Positive Sample
mit	Automobil	:= 1	
Haste	Automobil	:= 0	Negative Sample
oben	Automobil	:= 0	
auf	Automobil	:= 0	

# Word Embedding

## Training Skip-Gram Word Embedding with Negative Sampling

Training is computational expensive.

Ein Elektroauto ist ein Automobil mit elektrischem Antrieb.

$P(+|t, c)$

Positive Sample

Context	Target		
ein	Automobil	$:= 1$	} Positive Sample
mit	Automobil	$:= 1$	

# Word Embedding

## Training Skip-Gram Word Embedding with Negative Sampling

Training is computational expensive.

Ein Elektroauto ist ein Automobil mit elektrischem Antrieb.

This is a  
binary  
classifier {  $P(+|t, c)$   
Positive Sample

Context	Target		
ein	Automobil	$:= 1$	} Positive Sample
mit	Automobil	$:= 1$	

# Word Embedding

## Training Skip-Gram Word Embedding with Negative Sampling

Training is computational expensive. Replace softmax with sigmoid.

Ein Elektroauto ist ein Automobil mit elektrischem Antrieb.

This is a  
binary  
classifier  $\left\{ \begin{array}{l} P(+|t, c) \\ \text{Positive Sample} \end{array} \right.$



$$P(+|t, c) = \frac{1}{1 + e^{-t \cdot c}}$$

Context	Target		
ein	Automobil	$:= 1$	} Positive Sample
mit	Automobil	$:= 1$	

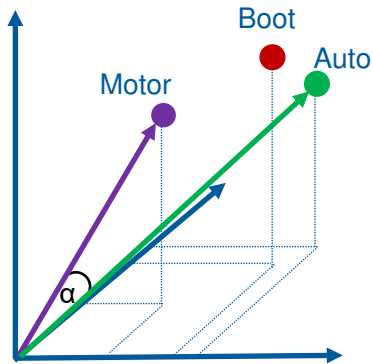
Sigmoid:

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

Training is computational expensive. Replace softmax with sigmoid.

Ein Elektroauto ist **ein** **Automobil** mit elektrischem Antrieb.

This is a  
binary  
classifier  $\left\{ \begin{array}{l} P(+|t, c) \\ \text{Positive Sample} \end{array} \right.$



$$P(+|t, c) = \frac{1}{1 + e^{-t \cdot c}}$$

$$\text{Similarity}(t, c) \approx t \cdot c$$

Context  
ein  
mit

Target  
Automobil  
Automobil

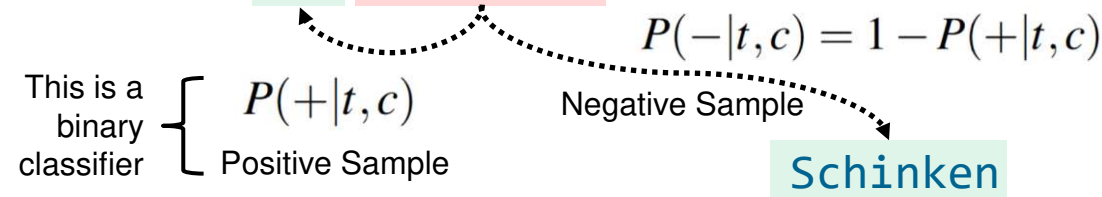
$\begin{array}{l} := 1 \\ := 1 \end{array} \left. \vphantom{\begin{array}{l} := 1 \\ := 1 \end{array}} \right\} \text{Positive Sample}$

Cosine Similarity:

$$\cos \angle(\vec{a}, \vec{b}) = \frac{\vec{a} \cdot \vec{b}}{|\vec{a}| |\vec{b}|}$$

Training is computational expensive. A (binary) classifier needs “negative samples”.

Ein Elektroauto ist ein Automobil mit elektrischem Antrieb.



$$P(+|t, c) = \frac{1}{1 + e^{-t \cdot c}}$$

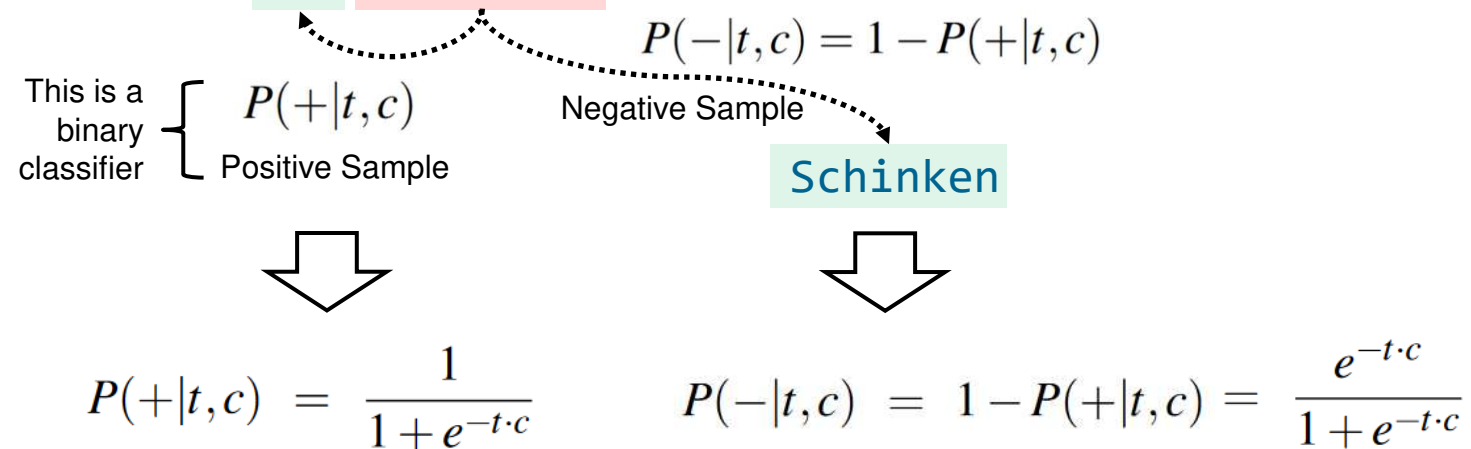
Context	Target		
ein	Automobil	:= 1	Positive Sample
mit	Automobil	:= 1	
Haste	Automobil	:= 0	Negative Sample
oben	Automobil	:= 0	
auf	Automobil	:= 0	

Sigmoid:

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

Training is computational expensive. A (binary) classifier needs “negative samples”.

Ein Elektroauto ist ein Automobil mit elektrischem Antrieb.

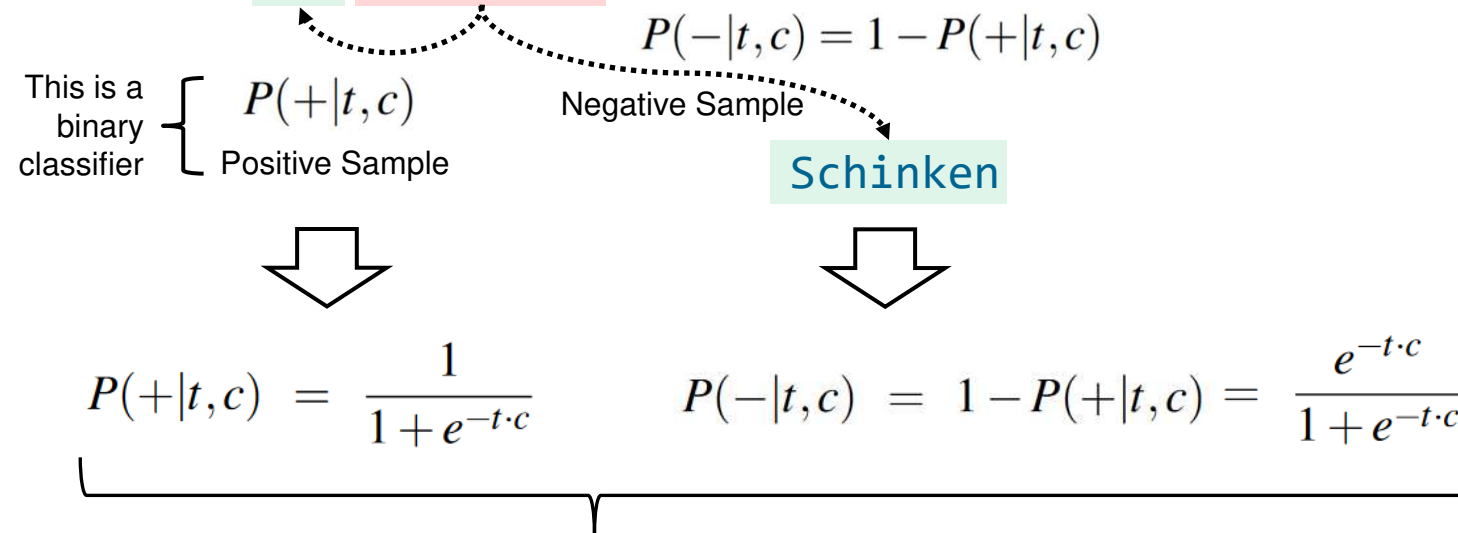


Sigmoid:

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

- Maximize the similarity of the target word, context word pairs (t,c) drawn from the positive examples
- Minimize the similarity of the (t,c) pairs drawn from the negative examples.

Ein Elektroauto ist **ein** **Automobil** mit elektrischem Antrieb.



$$L(\theta) = \log \frac{1}{1 + e^{-c \cdot t}} + \sum_{i=1}^k \log \frac{1}{1 + e^{n_i \cdot t}} \quad \text{For k noise samples}$$

Sigmoid:

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$



Training is computational expensive.

Ein Elektroauto ist ein Automobil mit elektrischem Antrieb.

Remember: We have a context window

$$P(+|t, c_{1:k}) = \prod_{i=1}^k \frac{1}{1 + e^{-t \cdot c_i}}$$

$$\log P(+|t, c_{1:k}) = \sum_{i=1}^k \log \frac{1}{1 + e^{-t \cdot c_i}}$$

Independents of context words assumed

This is a  
binary  
classifier

$P(+|t, c)$   
Positive Sample

Negative Sample

Schinken

$$P(+|t, c) = \frac{1}{1 + e^{-t \cdot c}}$$

$$P(-|t, c) = 1 - P(+|t, c) = \frac{e^{-t \cdot c}}{1 + e^{-t \cdot c}}$$

$$L(\theta) = \log \frac{1}{1 + e^{-c \cdot t}} + \sum_{i=1}^k \log \frac{1}{1 + e^{n_i \cdot t}} \quad \text{For } k \text{ noise samples}$$

Sigmoid:

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

Training is computational expensive.

Ein Elektroauto ist ein Automobil mit elektrischem Antrieb.

Remember: We have a context window

$$P(+|t, c_{1:k}) = \prod_{i=1}^k \frac{1}{1 + e^{-t \cdot c_i}}$$

$$\log P(+|t, c_{1:k}) = \sum_{i=1}^k \log \frac{1}{1 + e^{-t \cdot c_i}}$$

Independents of context words assumed

This is a  
binary  
classifier

$P(+|t, c)$   
Positive Sample

$$P(-|t, c) = 1 - P(+|t, c)$$

Negative Sample

Schinken

How to select the negative  
samples from the vocabulary?

$$P(w) = \frac{\text{count}(w)}{\sum_{w'} \text{count}(w')}$$

Uni-Gram Probabilities

$$P(+|t, c) = \frac{1}{1 + e^{-t \cdot c}}$$

$$P(-|t, c) = 1 - P(+|t, c) = \frac{e^{-t \cdot c}}{1 + e^{-t \cdot c}}$$

$$L(\theta) = \log \frac{1}{1 + e^{-c \cdot t}} + \sum_{i=1}^k \log \frac{1}{1 + e^{n_i \cdot t}} \quad \text{For } k \text{ noise samples}$$

Sigmoid:

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

Training is computational expensive.

Ein Elektroauto ist ein Automobil mit elektrischem Antrieb.

Remember: We have a context window

$$P(+|t, c_{1:k}) = \prod_{i=1}^k \frac{1}{1 + e^{-t \cdot c_i}}$$

$$\log P(+|t, c_{1:k}) = \sum_{i=1}^k \log \frac{1}{1 + e^{-t \cdot c_i}}$$

Independents of context words assumed

This is a  
binary  
classifier

$P(+|t, c)$   
Positive Sample

$$P(-|t, c) = 1 - P(+|t, c)$$

Negative Sample

Schinken

How to select the negative  
samples from the vocabulary?

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

Weighted Uni-Gram Probabilities  
Rare words:  $P_{\alpha}(w) > P(w)$

$$P(+|t, c) = \frac{1}{1 + e^{-t \cdot c}}$$

$$P(-|t, c) = 1 - P(+|t, c) = \frac{e^{-t \cdot c}}{1 + e^{-t \cdot c}}$$

$$L(\theta) = \log \frac{1}{1 + e^{-c \cdot t}} + \sum_{i=1}^k \log \frac{1}{1 + e^{n_i \cdot t}} \quad \text{For } k \text{ noise samples}$$

Sigmoid:

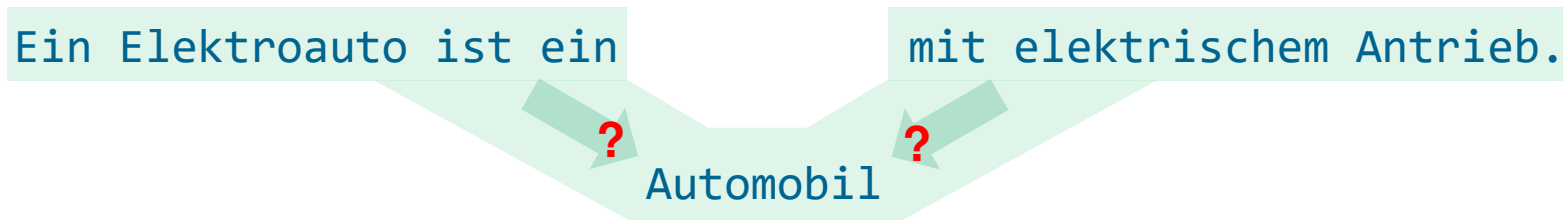
$$\sigma(x) = \frac{1}{1 + e^{-x}}$$



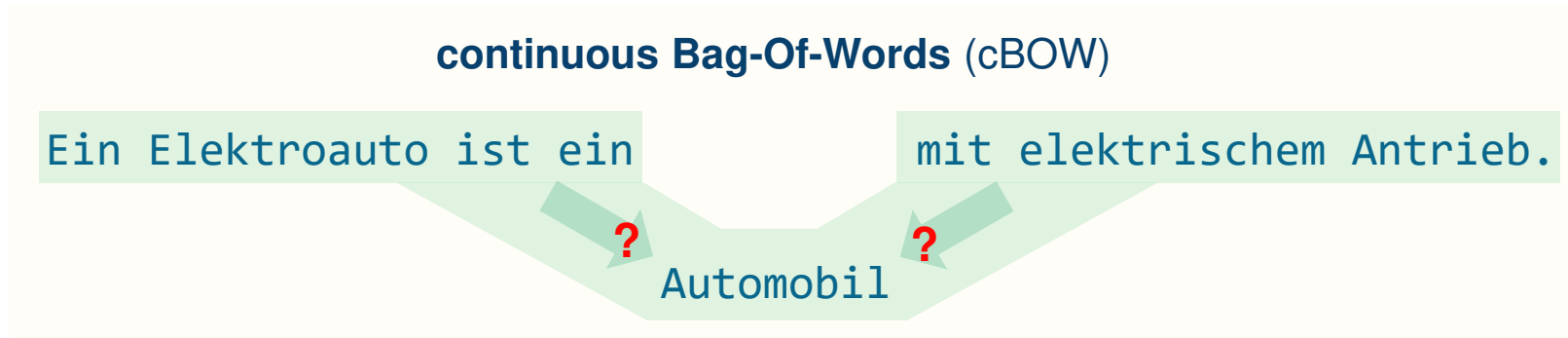
Word embedding is any of a set of language modeling and feature learning techniques in natural language processing where words from the vocabulary are mapped to vectors of real numbers.

Ein Elektroauto ist ein Automobil mit elektrischem Antrieb.

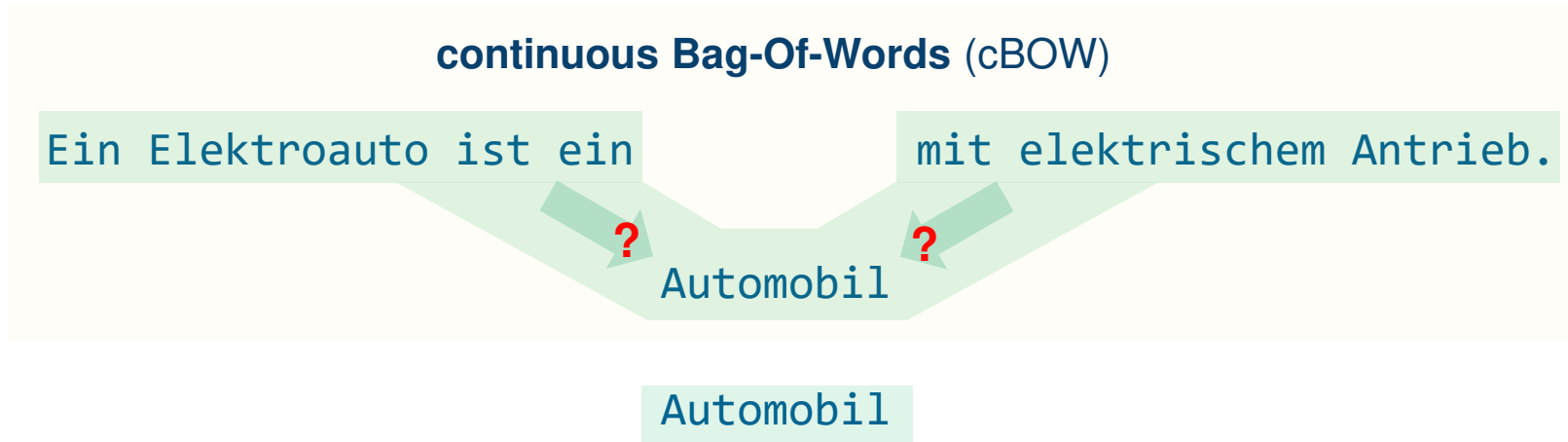
Word embedding is any of a set of language modeling and feature learning techniques in natural language processing where words from the vocabulary are mapped to vectors of real numbers.



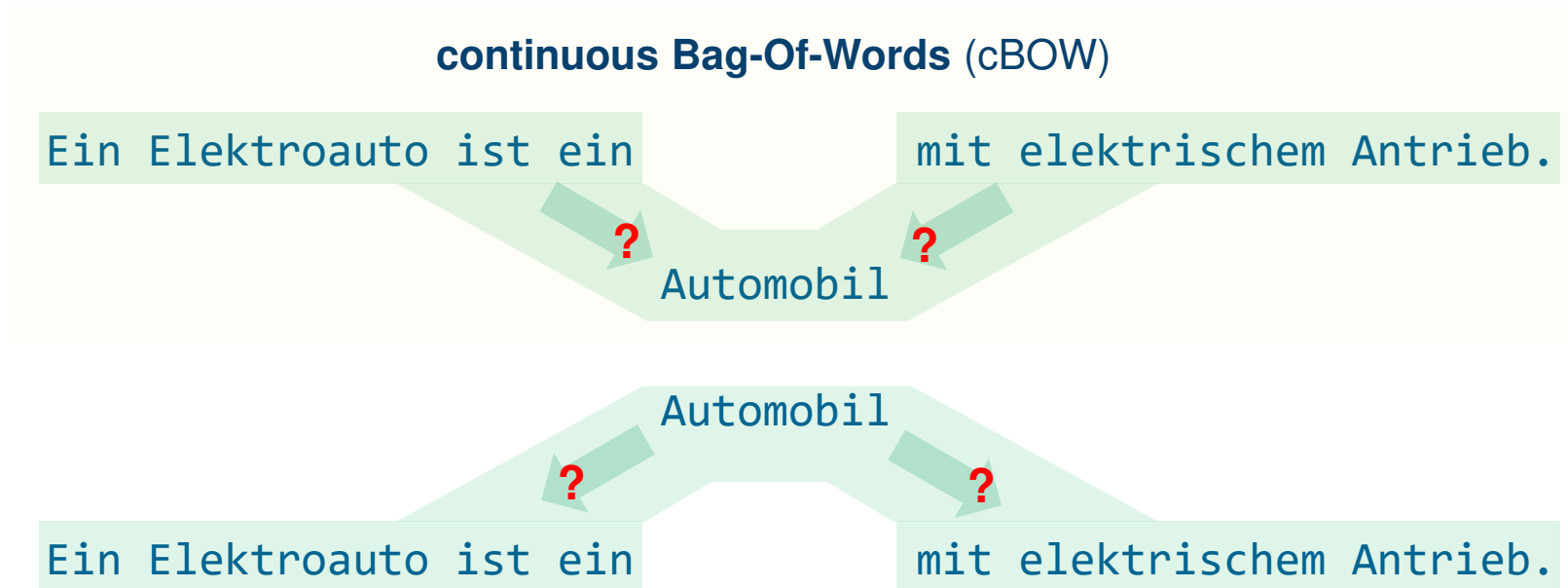
Word embedding is any of a set of language modeling and feature learning techniques in natural language processing where words from the vocabulary are mapped to vectors of real numbers.



Word embedding is any of a set of language modeling and feature learning techniques in natural language processing where words from the vocabulary are mapped to vectors of real numbers.

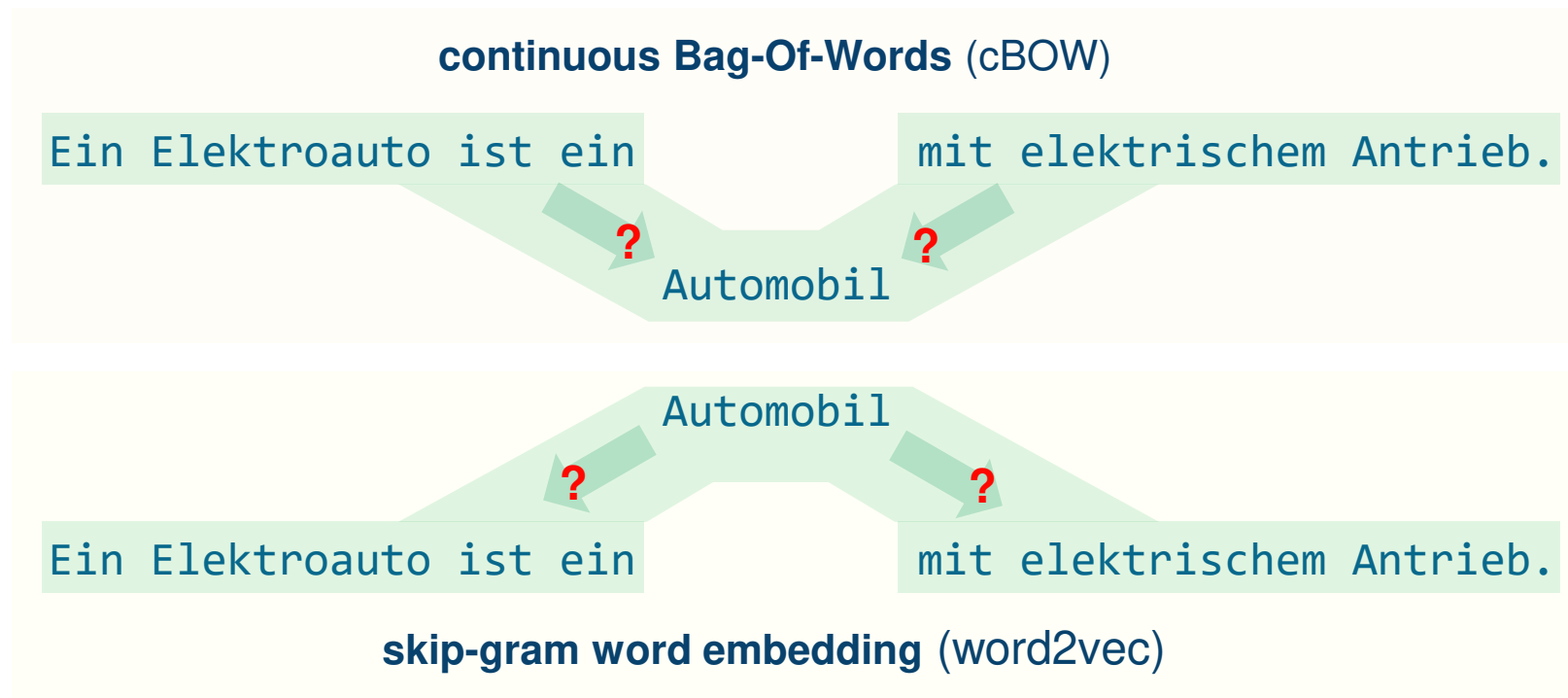


Word embedding is any of a set of language modeling and feature learning techniques in natural language processing where words from the vocabulary are mapped to vectors of real numbers.





Word embedding is any of a set of language modeling and feature learning techniques in natural language processing where words from the vocabulary are mapped to vectors of real numbers.





The input is always one word of the sentence

Ein Elektroauto ist ein **Automobil** mit elektrischem Antrieb.

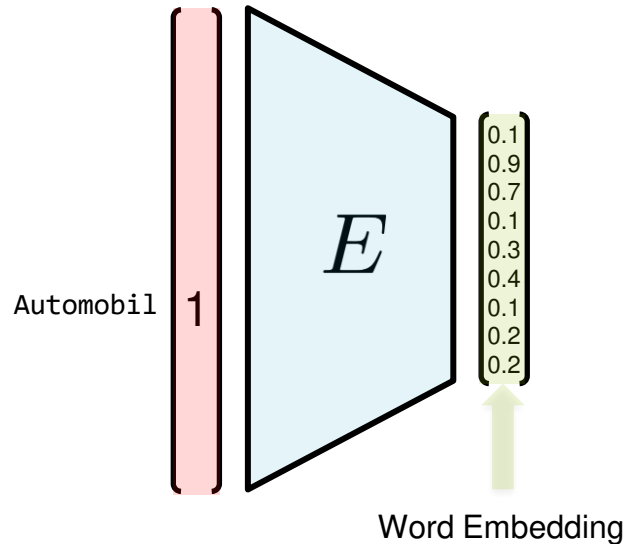
for **targetword** in sentence:

Automobil 1

$$x_i = Ew_i$$

Compute the word embedding of given word

Ein Elektroauto ist ein **Automobil** mit elektrischem Antrieb.



for **targetword** in sentence:

```
embedding = matmul(E, targetword)
```

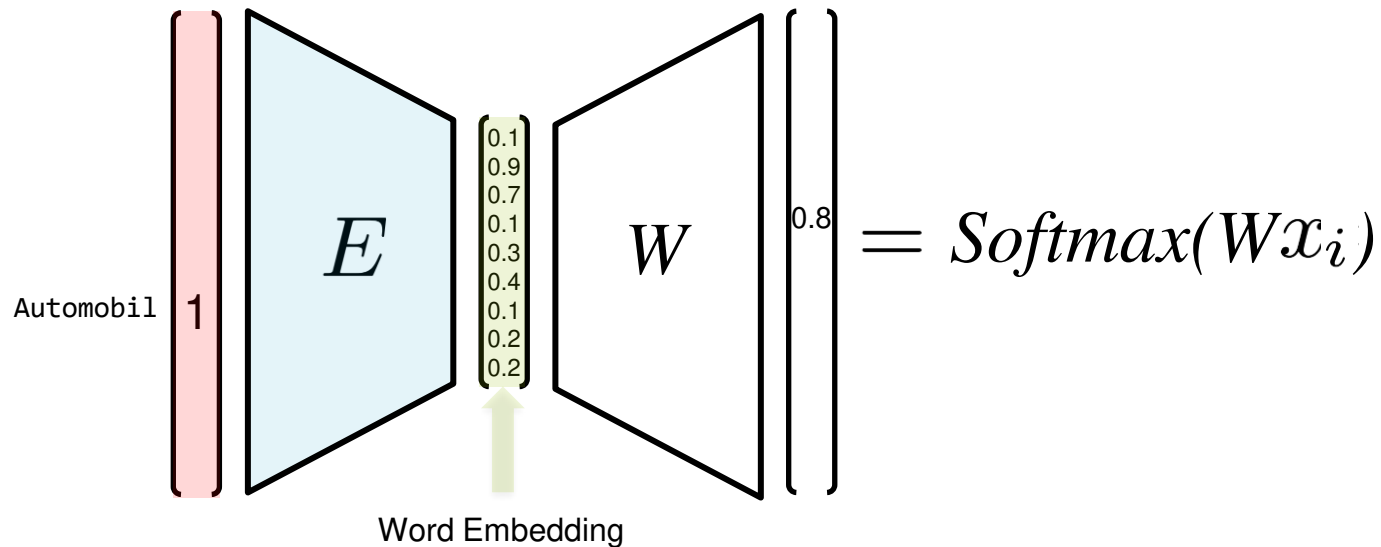
$$x_i = E w_i$$

# Word Embedding

## Skip-Gram Word Embedding

Make a prediction using the word embedding aka “use the word embedding”

Ein Elektroauto ist ein **Automobil** mit elektrischem Antrieb.



for **targetword** in sentence:

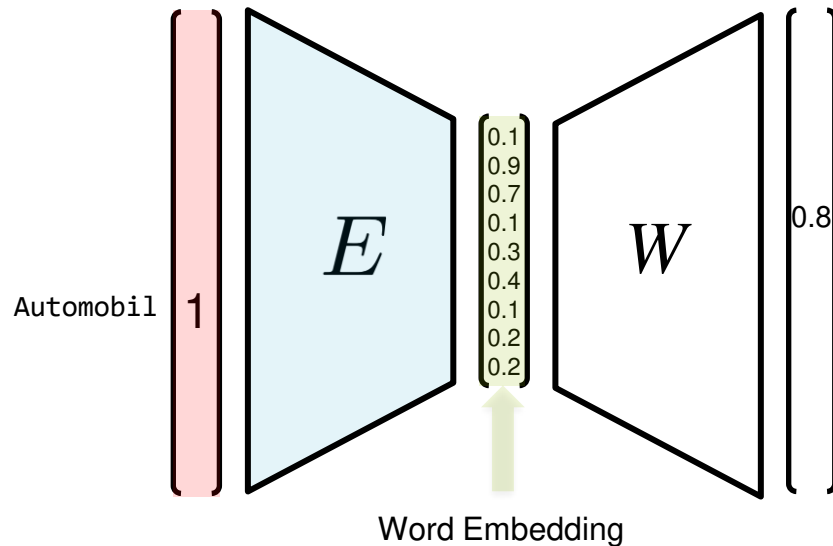
```
embedding = matmul(E, targetword)
tmp = matmul(W, embedding)
predicted_contextword = softmax(tmp)
```

$$x_i = Ew_i$$

One possible prediction: Predict the context of the word

Context: +/- 2 Words

Ein Elektroauto ist ein **Automobil** mit elektrischem Antrieb.



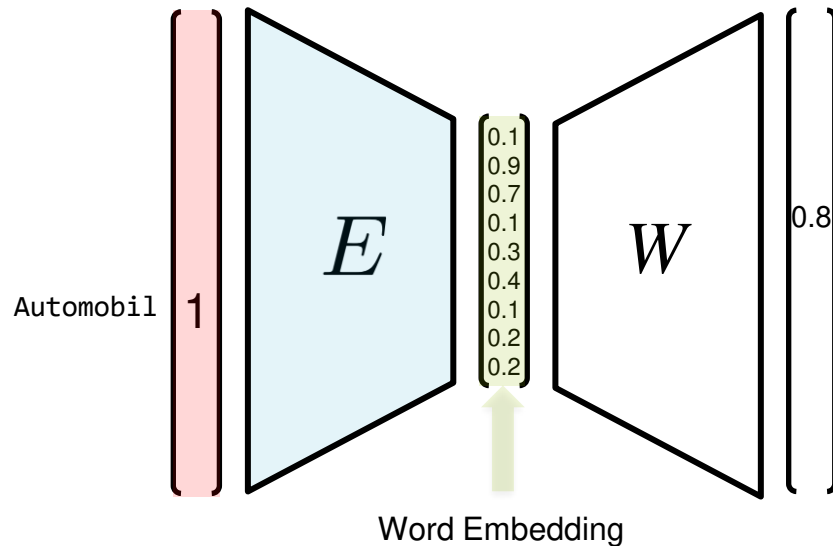
```
for targetword in sentence:  
    for contextword arround targetword:  
        embedding = matmul(E,targetword)  
        tmp = matmul(W, embedding)  
        predicted_contextword = softmax(tmp)
```

$$x_i = Ew_i$$

One possible prediction: Predict the context of the word

Context: +/- 2 Words

Ein Elektroauto ist ein **Automobil** mit **elektrischem** Antrieb.



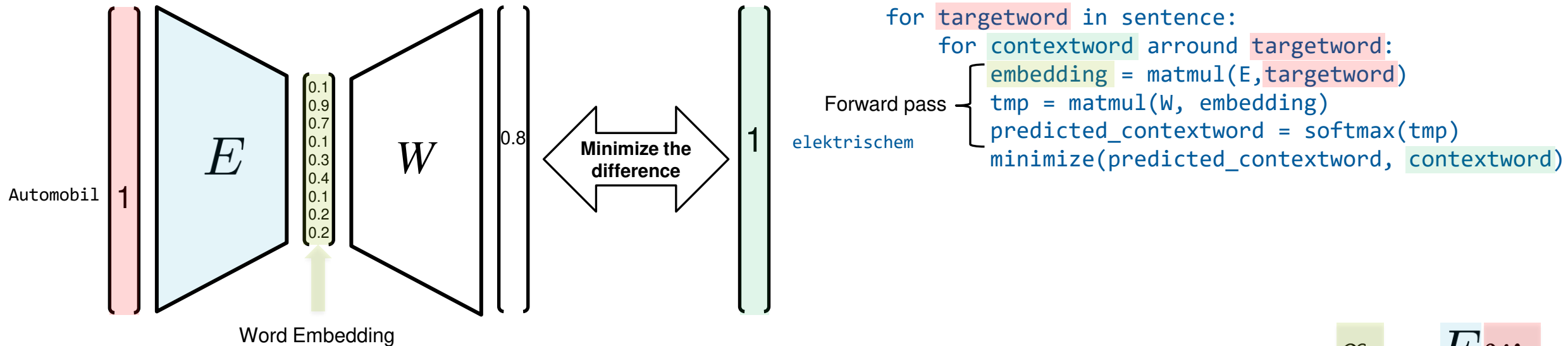
```
for targetword in sentence:  
    for contextword around targetword:  
        embedding = matmul(E, targetword)  
        tmp = matmul(W, embedding)  
        predicted_contextword = softmax(tmp)
```

Forward pass {  
elektrischem

$$x_i = Ew_i$$

Minimize the error between prediction of the context and the real context by updating the network

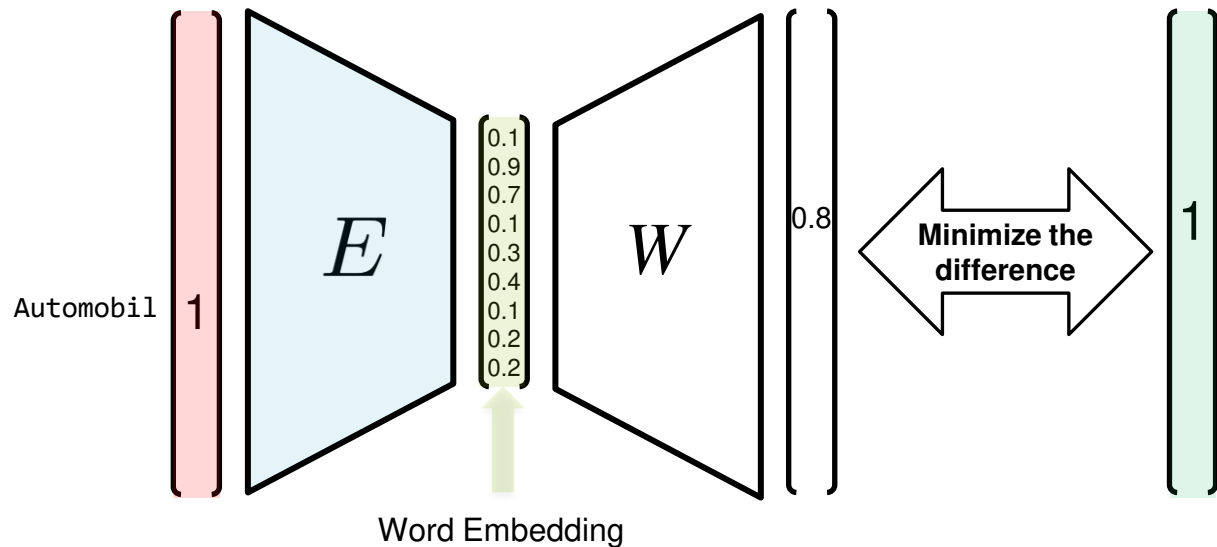
Ein Elektroauto ist ein **Automobil** mit **elektrischem** Antrieb.



$$x_i = E w_i$$

Minimize the error between prediction of the context and the real context by updating the network

Ein Elektroauto ist ein **Automobil** mit elektrischem Antrieb.



```
for targetword in sentence:
    for contextword around targetword:
        embedding = matmul(E, targetword)
        tmp = matmul(W, embedding)
        predicted_contextword = softmax(tmp)
        minimize(predicted_contextword, contextword)
```

Forward pass {

Backward pass {

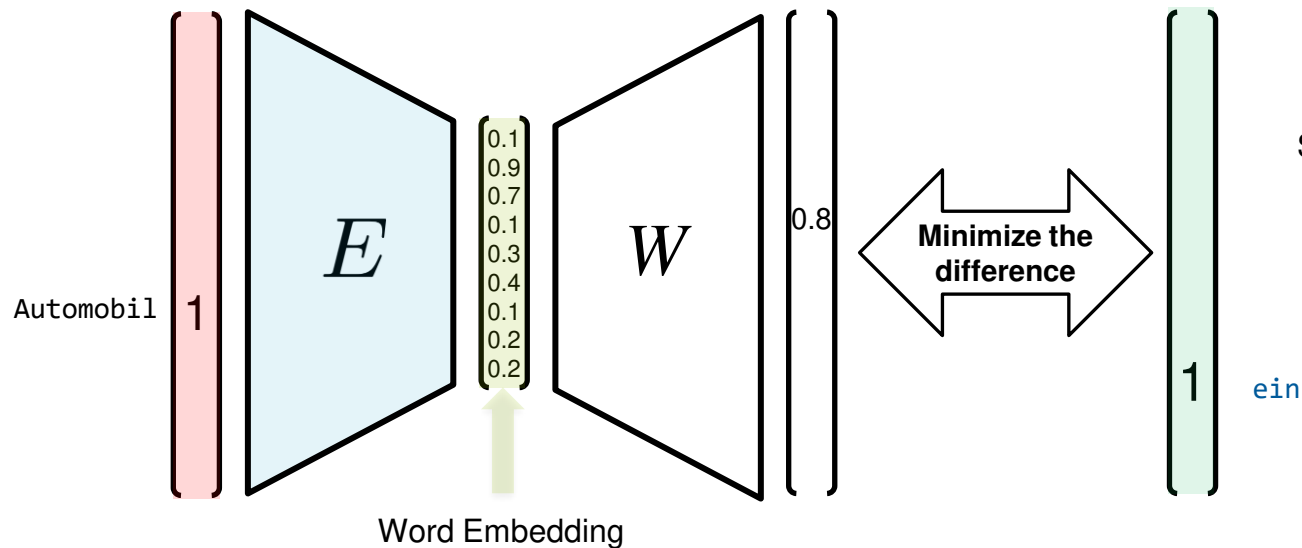
Updates {

$$x_i = E w_i$$



Considering the complete context of a word is almost impossible: sample.

Ein Elektroauto ist ein Automobil mit elektrischem Antrieb.



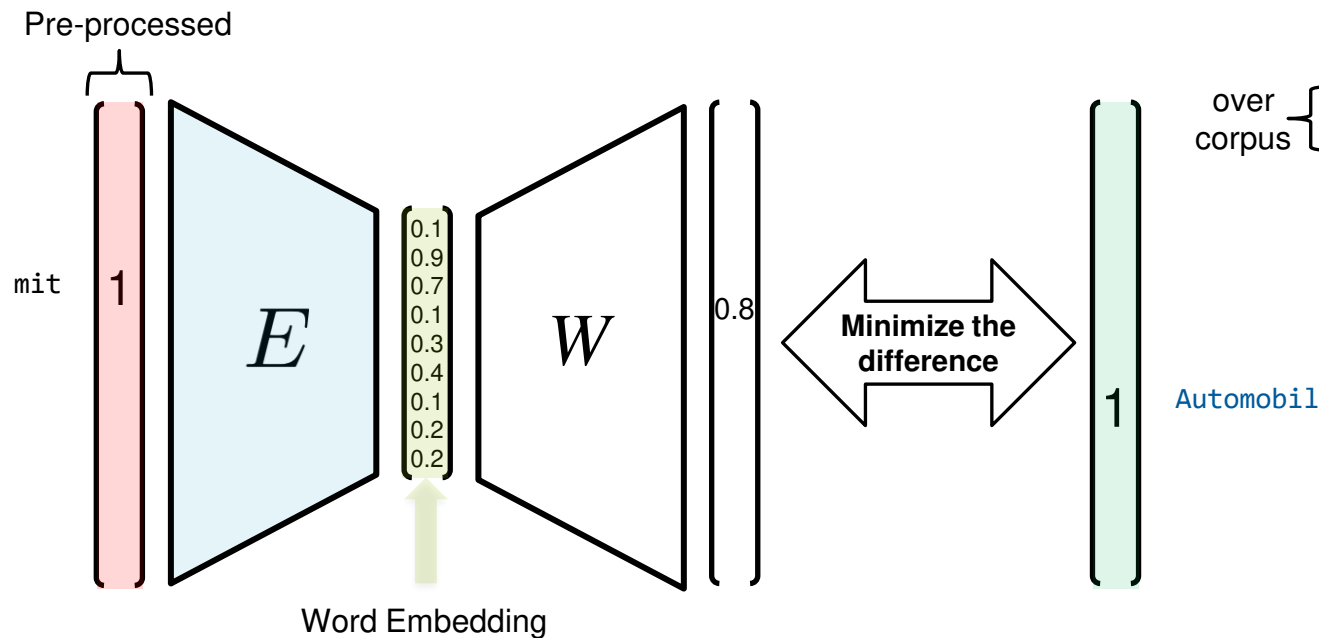
Sampling required

```
for targetword in sentence:
    for contextword around targetword:
        embedding = matmul(E, targetword)
        tmp = matmul(W, embedding)
        predicted_contextword = softmax(tmp)
        minimize(predicted_contextword, contextword)
```

$$x_i = E w_i$$

Repeat the process for any word in the corpus

Ein Elektroauto ist ein Automobil mit elektrischem Antrieb.



over corpus {

```
for targetword in sentence:
    for contextword around targetword:
        embedding = matmul(E, targetword)
        tmp = matmul(W, embedding)
        predicted_contextword = softmax(tmp)
        minimize(predicted_contextword, contextword)
```

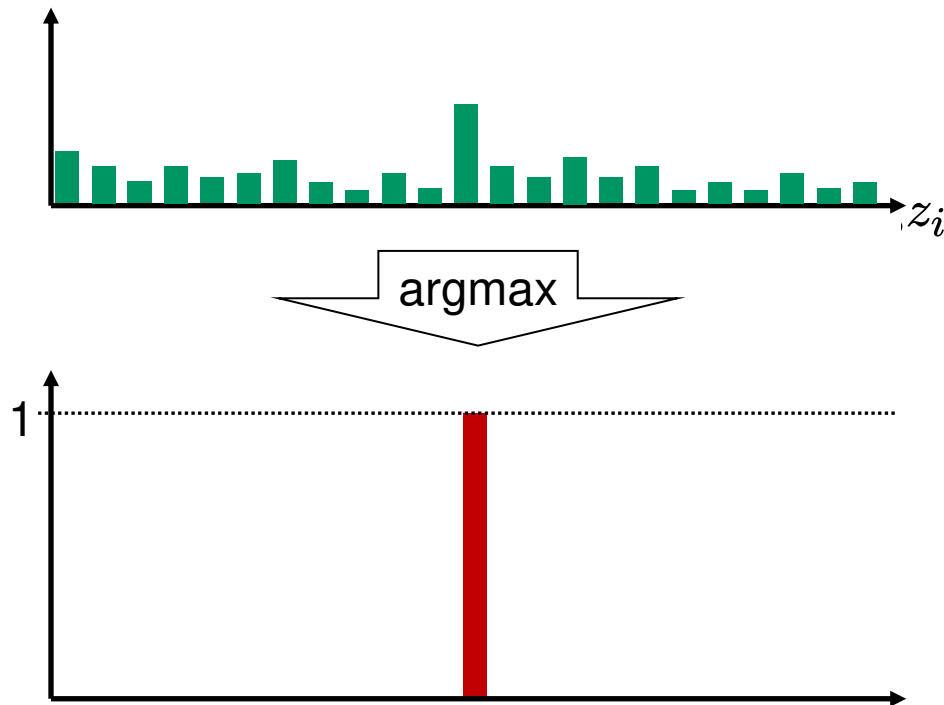
$$x_i = E w_i$$

# Softmax

max and argmax function



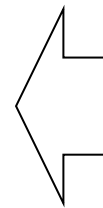
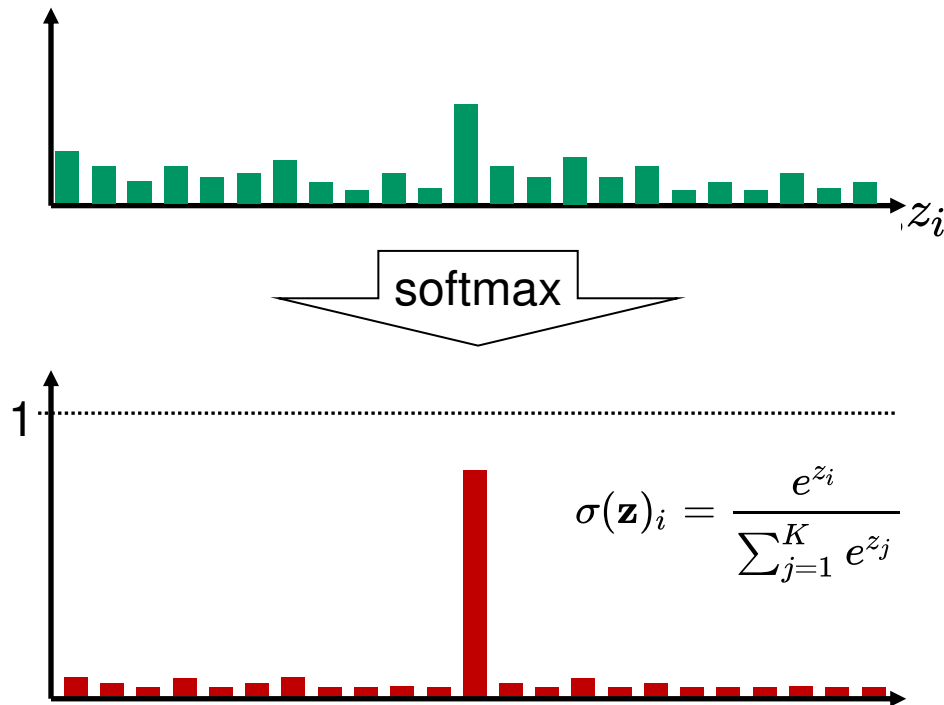
The softmax function is a function that turns a vector of  $K$  real values into a vector of  $K$  real values that sum to 1.



Where is the maximum?

All values are 0 except the one where the maximum is. This value is 1. Hence, the sum of all entries is 1, too.

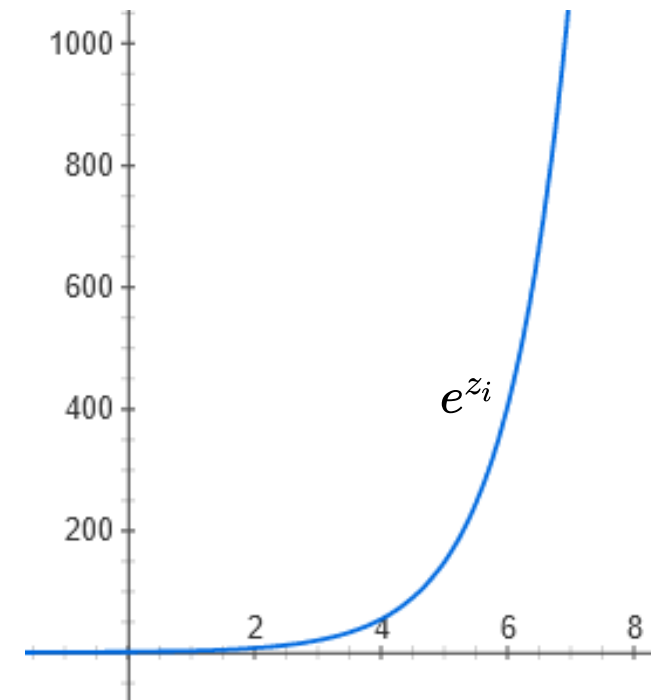
The softmax function is a function that turns a vector of  $K$  real values into a vector of  $K$  real values that sum to 1.



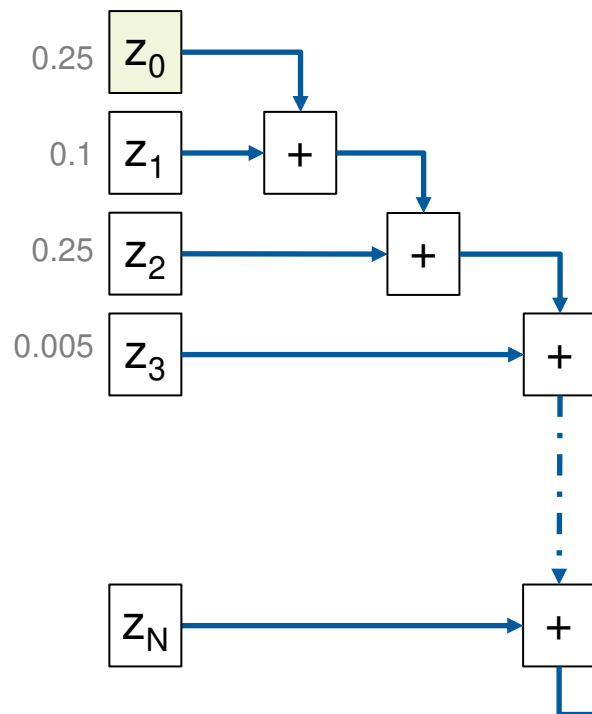
Where is the maximum?



Maximum is significantly larger than all other values. All other values are very close to 0.

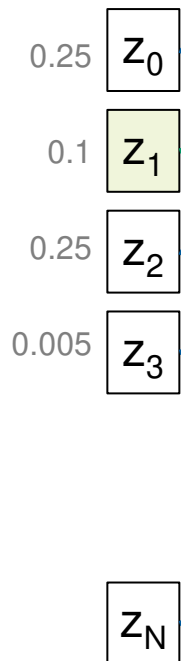


The softmax function is a function that turns a vector of  $K$  real values into a vector of  $K$  real values that sum to 1.

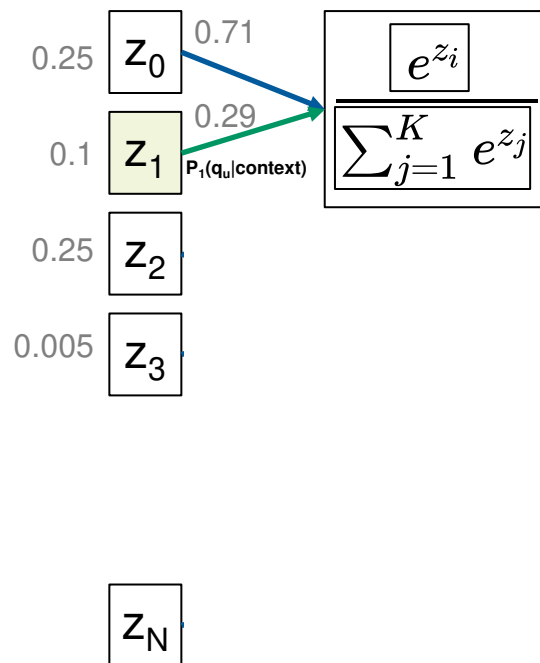


$$\sigma(\mathbf{z})_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}} \text{ for } i = 1, \dots, K \text{ and } \mathbf{z} = (z_1, \dots, z_K) \in \mathbb{R}^K$$

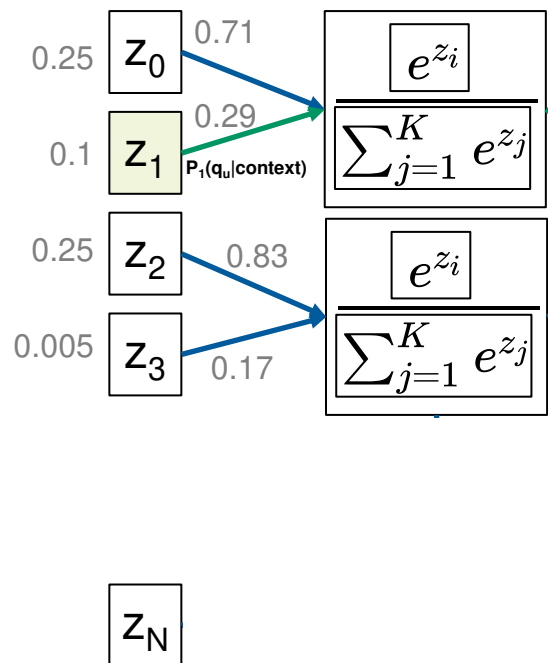
An approximation inspired by binary trees that replaces the flat softmax layer with a hierarchical layer that has the words as leaves.



An approximation inspired by binary trees that replaces the flat softmax layer with a hierarchical layer that has the words as leaves.

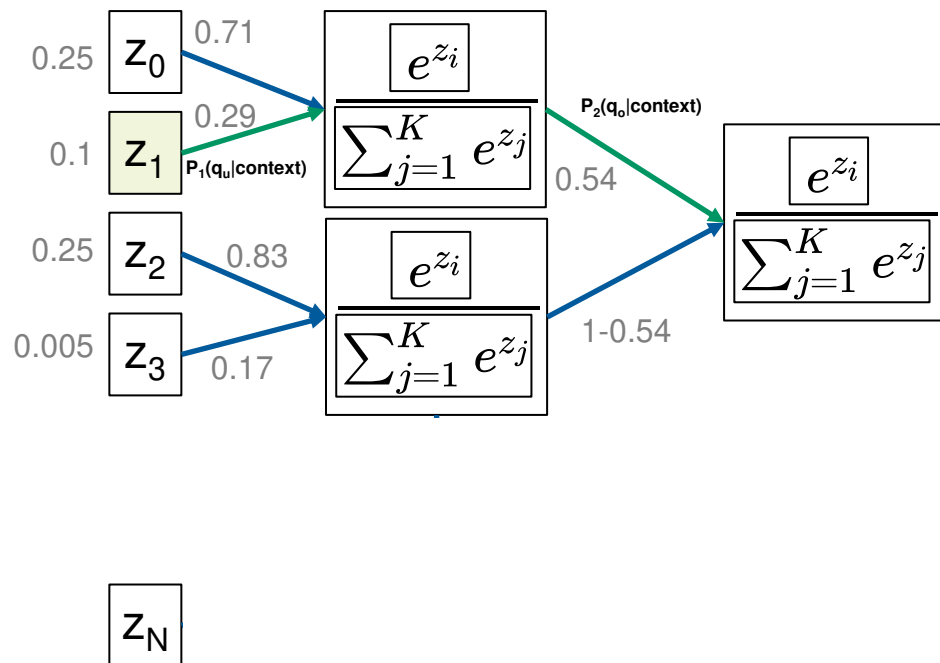


An approximation inspired by binary trees that replaces the flat softmax layer with a hierarchical layer that has the words as leaves.

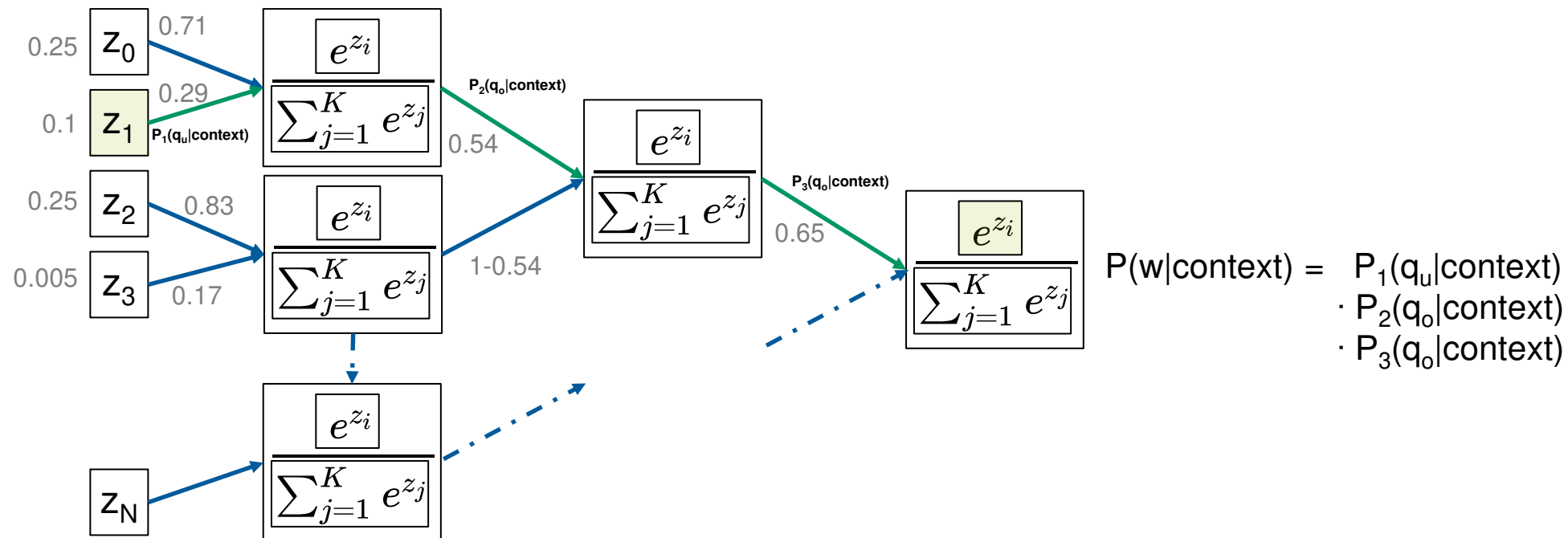




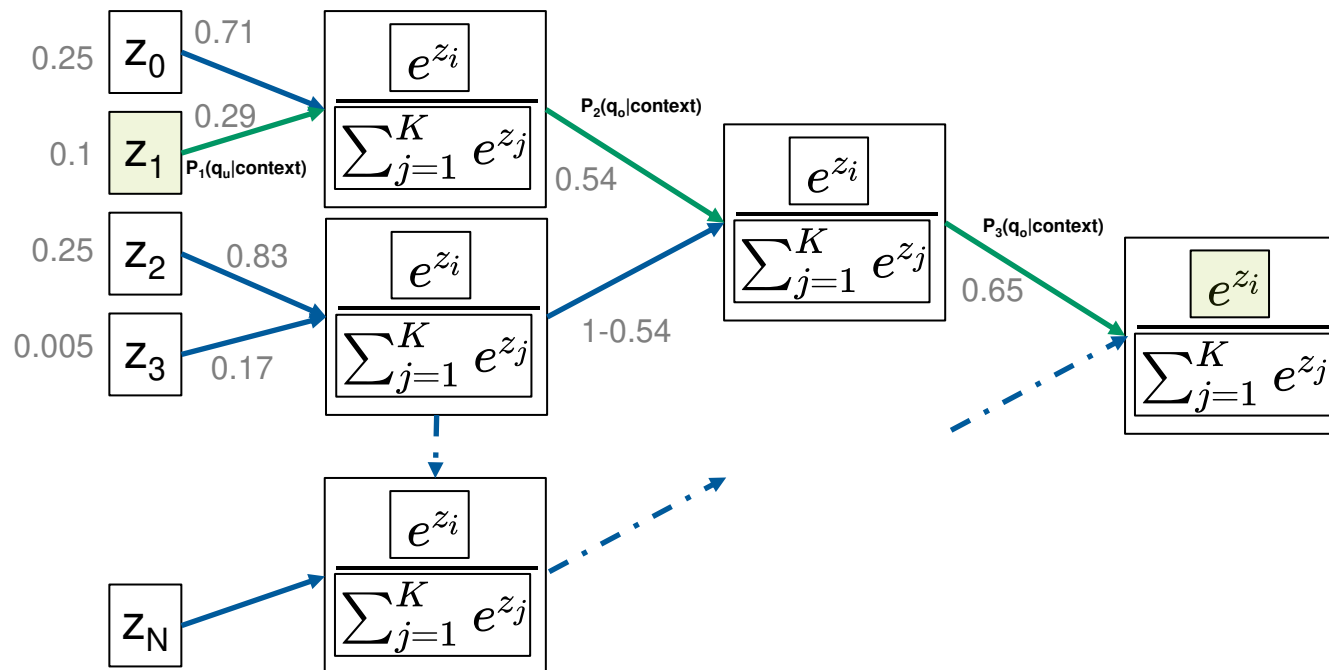
An approximation inspired by binary trees that replaces the flat softmax layer with a hierarchical layer that has the words as leaves.



An approximation inspired by binary trees that replaces the flat softmax layer with a hierarchical layer that has the words as leaves.



An approximation inspired by binary trees that replaces the flat softmax layer with a hierarchical layer that has the words as leaves.

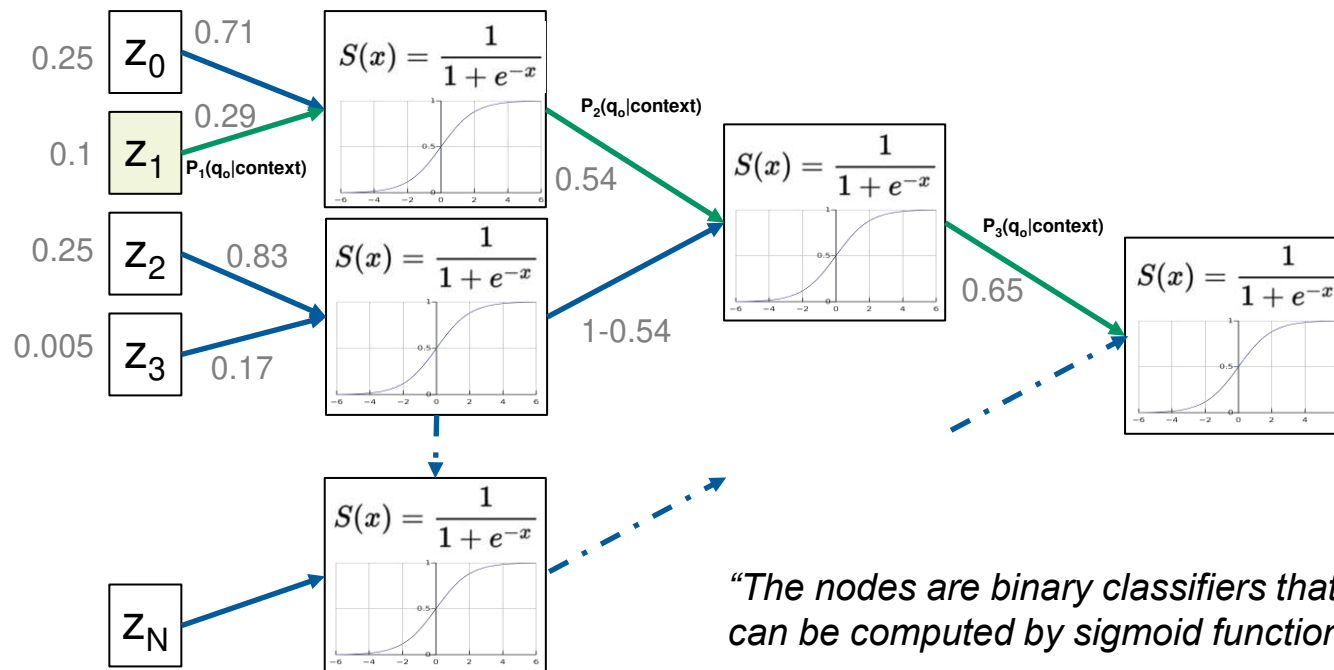


*“The probability of a word  $w$  given a context, is encoded as path from root to word  $w$ .”*

Complexity reduces to  $\log(|V|)$  for binary trees

$$P(w|\text{context}) = P_1(q_u|\text{context}) \cdot P_2(q_o|\text{context}) \cdot P_3(q_o|\text{context})$$

An approximation inspired by binary trees that replaces the flat softmax layer with a hierarchical layer that has the words as leaves.



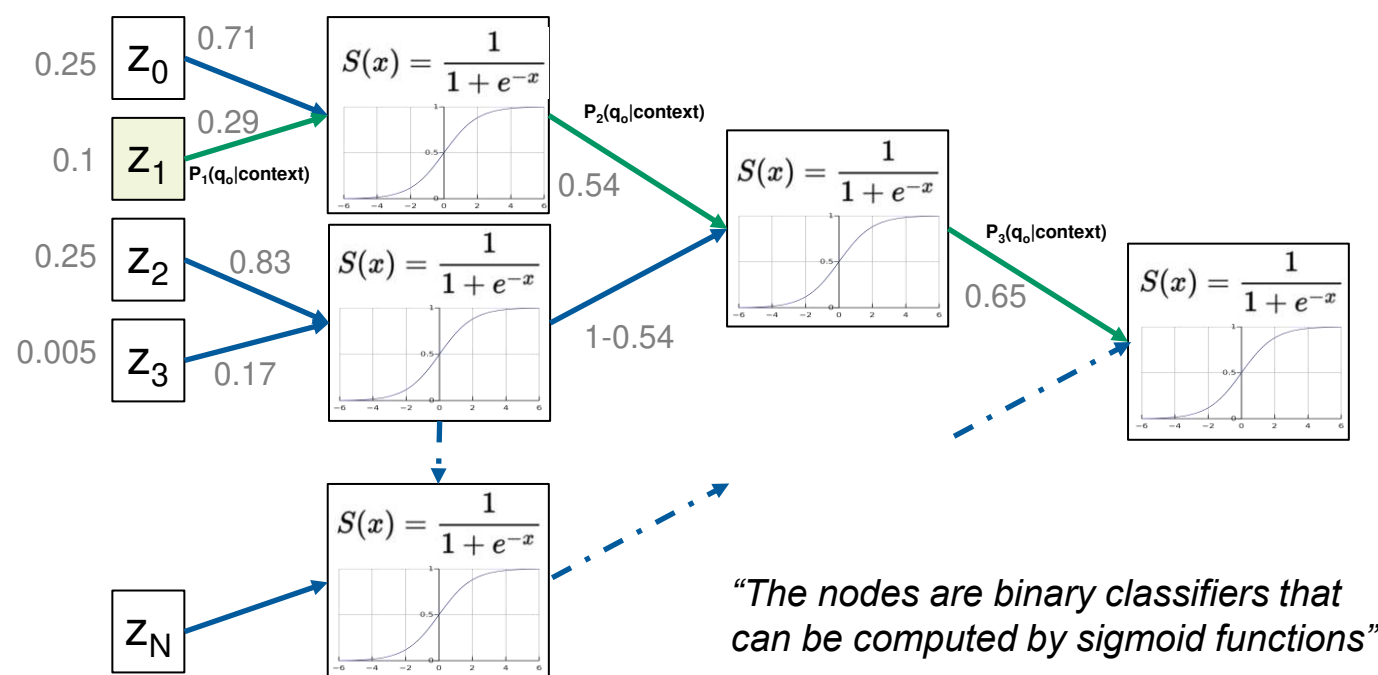
*"The probability of a word  $w$  given a context, is encoded as path from root to word  $w$ ."*

Complexity reduces to  $\log(|V|)$  for binary trees

$$P(w|\text{context}) = P_1(q_u|\text{context}) \cdot P_2(q_o|\text{context}) \cdot P_3(q_o|\text{context})$$

*"The nodes are binary classifiers that can be computed by sigmoid functions"*

An approximation inspired by binary trees that replaces the flat softmax layer with a hierarchical layer that has the words as leaves.



*"The probability of a word  $w$  given a context, is encoded as path from root to word  $w$ ."*

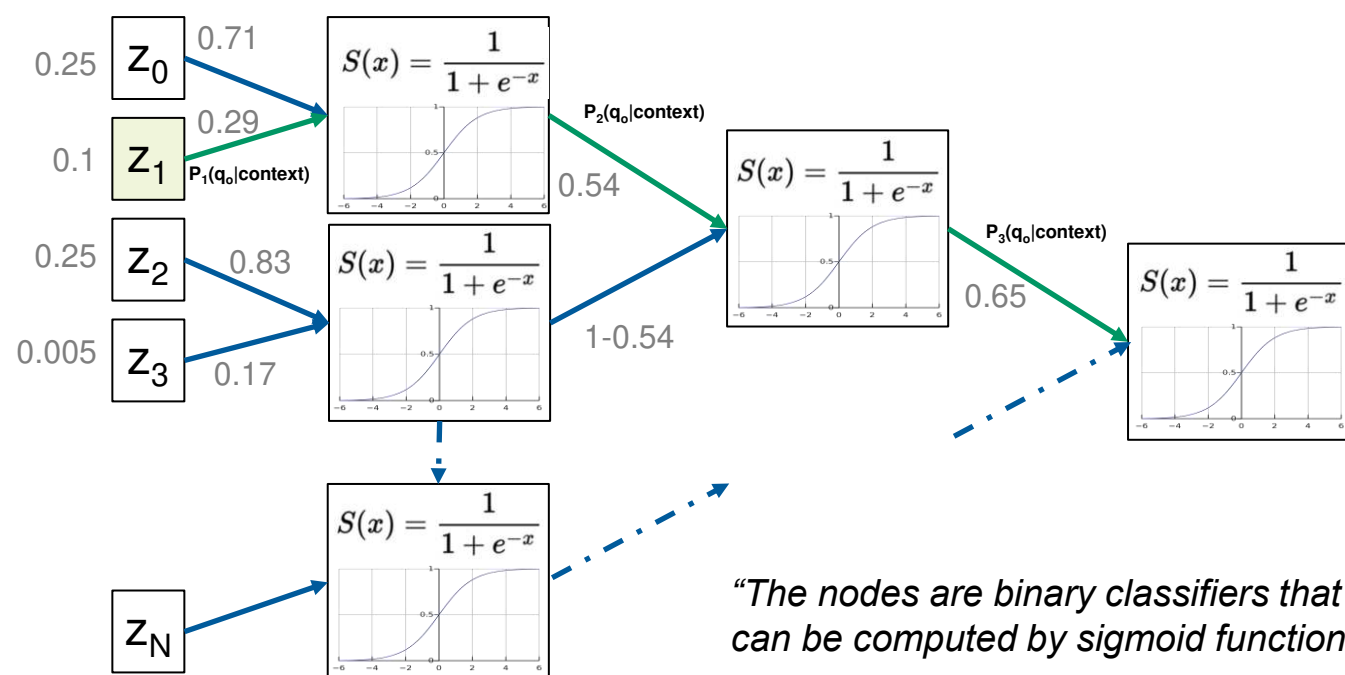
$$P(w|\text{context}) = P_1(q_u|\text{context}) \cdot P_2(q_o|\text{context}) \cdot P_3(q_o|\text{context})$$

Complexity reduces to  $\log(|V|)$  for binary trees

What about none binary trees?

*"The nodes are binary classifiers that can be computed by sigmoid functions"*

An approximation inspired by binary trees that replaces the flat softmax layer with a hierarchical layer that has the words as leaves.



*"The probability of a word  $w$  given a context, is encoded as path from root to word  $w$ ."*

$$P(w|\text{context}) = P_1(q_u|\text{context}) \cdot P_2(q_o|\text{context}) \cdot P_3(q_o|\text{context})$$

*"The nodes are binary classifiers that can be computed by sigmoid functions"*

Complexity reduces to  $\log(|V|)$  for binary trees

What about none binary trees?

Still in research, e.g., wordnet is in discussion...

# Word Embedding

*t-distributed stochastic neighbor embedding (T-SNE)*



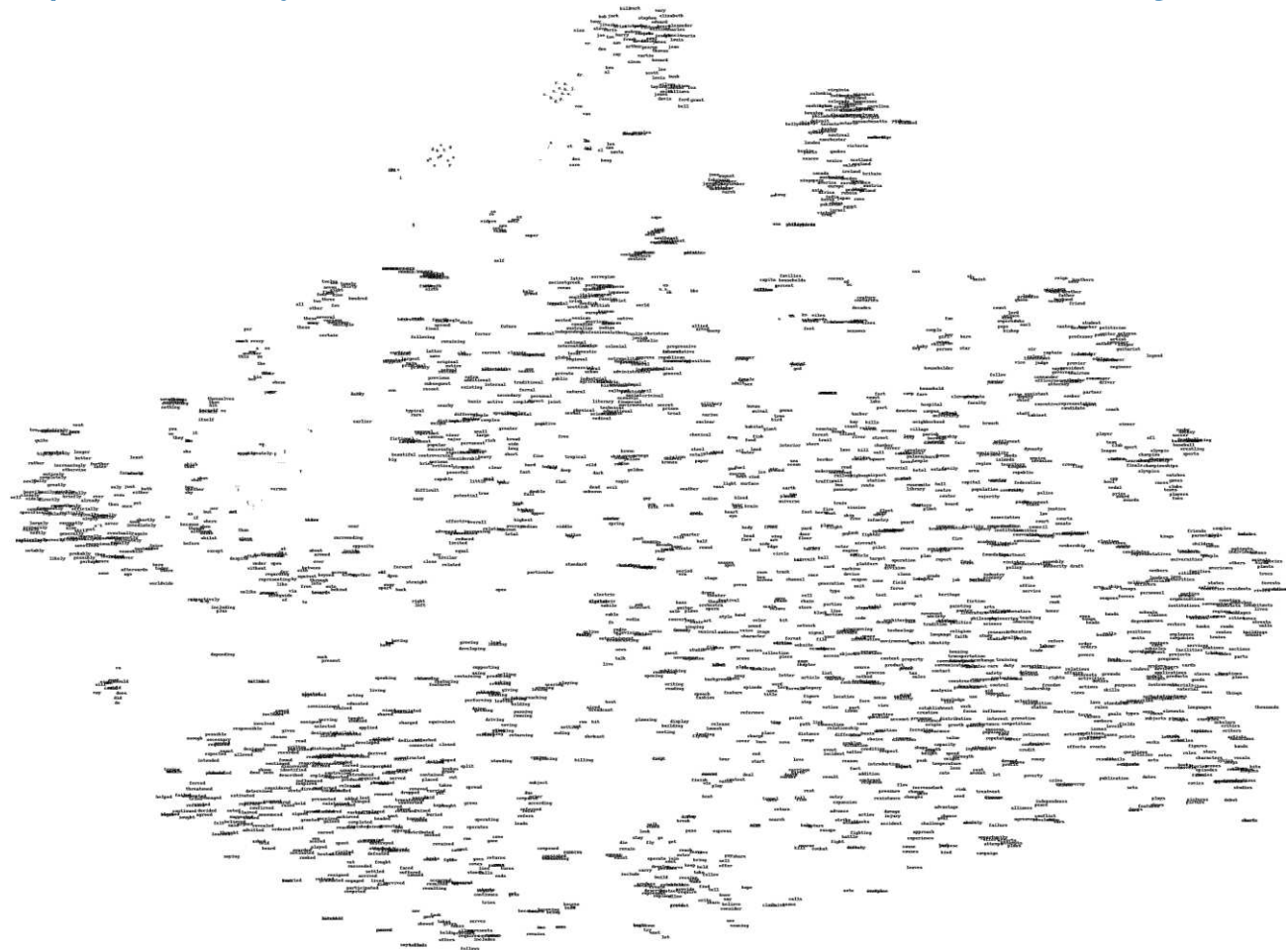
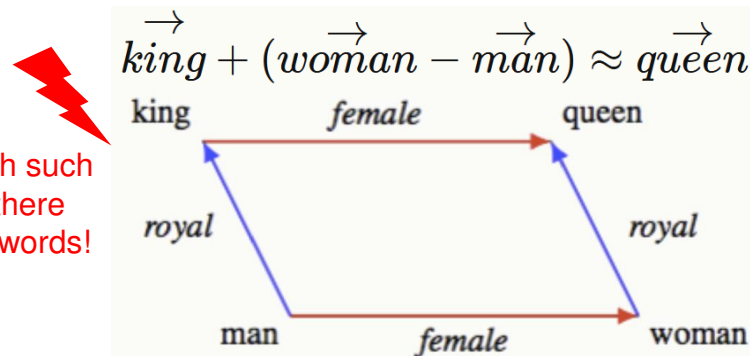
A technique for dimensionality reduction that is particularly well suited for the visualization of high-dimensional datasets.

Online version: <https://projector.tensorflow.org/>

Examples: <https://lvdmaaten.github.io/tsne/>

How To: <https://distill.pub/2016/misread-tsne/>

Be careful with such conclusions, there are million of words!



# Word Embedding

## Proposals

Embedding from a space with many dimensions per word to a continuous vector space with a much lower dimension. Generating a word embedding include neural networks, dimensionality reduction on word co-occurrence matrices, probabilistic models, explainable knowledge base method, and explicit representation in terms of the context in which words appear.

### (Classic) Word Embedding:

- Word2Vec: cBOW, Skip-Gram
- SGNS (skipgram with negative sampling) can be formalizes as:  $\text{SGNS} : \langle \vec{x}, \vec{y}_c \rangle = \text{PMI}(x, y) - \log k$
- GloVe (:= Global Vectors)
  - trained from global word-word co-occurrence statistics
  - <https://nlp.stanford.edu/projects/glove/>

$$\text{GloVe} : \langle \vec{x}, \vec{y}_c \rangle = \log X_{x,y} - b_x - b_y$$





Embedding from a space with many dimensions per word to a continuous vector space with a much lower dimension.

Generating a word embedding include neural networks, dimensionality reduction on word co-occurrence matrices, probabilistic models, explainable knowledge base method, and explicit representation in terms of the context in which words appear.

### (Classic) Word Embedding:

- Word2Vec: cBOW, Skip-Gram
- SGNS (skipgram with negative sampling) can be formalizes as:  $\text{SGNS} : \langle \vec{x}, \vec{y}_c \rangle = \text{PMI}(x, y) - \log k$
- GloVe (:= Global Vectors)  $\text{GloVe} : \langle \vec{x}, \vec{y}_c \rangle = \log X_{x,y} - b_x - b_y$ 
  - trained from global word-word co-occurrence statistics
  - <https://nlp.stanford.edu/projects/glove/>

### Context Sensitive Word Embedding (The representation for each word depends on the entire context in which it is used.)

- ELMo <https://allennlp.org/elmo>



Embedding from a space with many dimensions per word to a continuous vector space with a much lower dimension.

Generating a word embedding include neural networks, dimensionality reduction on word co-occurrence matrices, probabilistic models, explainable knowledge base method, and explicit representation in terms of the context in which words appear.

### (Classic) Word Embedding:

- Word2Vec: cBOW, Skip-Gram
- SGNS (skipgram with negative sampling) can be formalizes as:  $\text{SGNS} : \langle \vec{x}, \vec{y}_c \rangle = \text{PMI}(x, y) - \log k$
- GloVe (:= Global Vectors)  $\text{GloVe} : \langle \vec{x}, \vec{y}_c \rangle = \log X_{x,y} - b_x - b_y$ 
  - trained from global word-word co-occurrence statistics
  - <https://nlp.stanford.edu/projects/glove/>

### Context Sensitive Word Embedding (The representation for each word depends on the entire context in which it is used.)

- ELMo <https://allennlp.org/elmo>

### Non-Deterministic:

- T-SNE (t-distributed stochastic neighbor embedding)
  - [https://en.wikipedia.org/wiki/T-distributed\\_stochastic\\_neighbor\\_embedding](https://en.wikipedia.org/wiki/T-distributed_stochastic_neighbor_embedding) + see exercises