

PSY9511: Seminar 8

Sequence modelling (with an emphasis on language)

Esten H. Leonardsen

13.05.24



**UNIVERSITETET
I OSLO**

1. Introduction and motivation
2. Preprocessing
3. Bag of words
4. Vectorization
5. Recurrent neural networks
6. Transformers

Introduction



Introduction



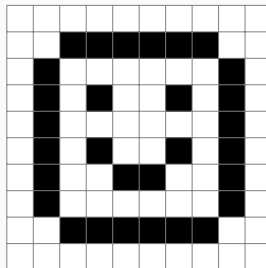
Age	Sex	Education	Salary
25	Male	12	40,000
30	Female	16	65,000
35	Male	14	55,000
40	Female	18	80,000
45	Male	16	75,000



Introduction

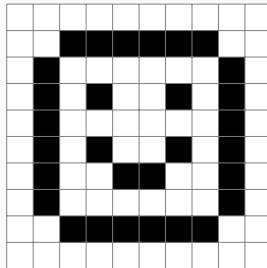


Age	Sex	Education	Salary
25	Male	12	40,000
30	Female	16	65,000
35	Male	14	55,000
40	Female	18	80,000
45	Male	16	75,000



The movie was great, the actors were awesome.

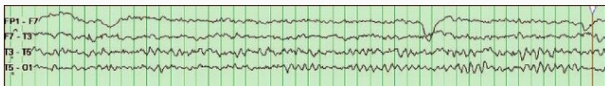
Age	Sex	Education	Salary
25	Male	12	40,000
30	Female	16	65,000
35	Male	14	55,000
40	Female	18	80,000
45	Male	16	75,000



The movie was great, the actors were awesome.

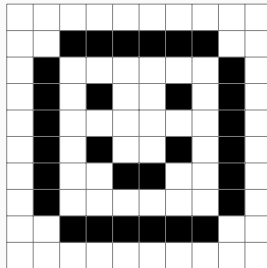


The movie was great, the actors were awesome.



The movie was great, the actors were awesome.

Age	Sex	Education	Salary
25	Male	12	40,000
30	Female	16	65,000
35	Male	14	55,000
40	Female	18	80,000
45	Male	16	75,000



Introduction



Age	Sex	Education	Salary
25	Male	12	40,000
30	Female	16	65,000
35	Male	14	55,000
40	Female	18	80,000
45	Male	16	75,000



Introduction



Age

35



Introduction

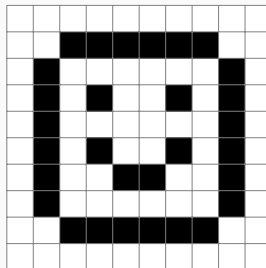


Age Sex

35 Male



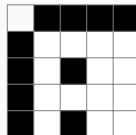
Introduction



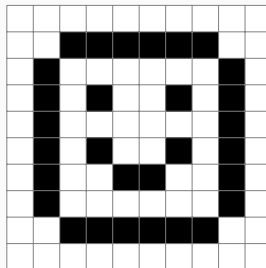
Introduction



Introduction



Introduction





The movie was great, the actors were awesome.



Introduction



movie





The movie was great, the actors were awesome.





The movie was great, the actors were awesome.

Positive

Negative



Introduction

The movie was great, the actors were awesome.

Positive

Negative

The movie was awful, the actors were horrible.



The movie was great, the actors were awesome.

La película fue genial, los actores fueron increíbles.





The movie was great, the actors were awesome.

La película fue genial, las actores fueron increíbles.





The movie was great, the actors were _____ .



Introduction



Introduction



The movie was _____ , the actors were _____ .





The movie was great, the actors were _____ .






The movie was great, the actors were awesome.



Introduction

The movie was great, the actors were awesome.



A diagram illustrating a dependency arc in a sentence. The sentence is "The movie was great, the actors were awesome." The words "great" and "awesome" are highlighted with black boxes. A curved arrow points from the box around "great" to the box around "awesome", indicating a dependency relationship between the two adjectives.





The movie was great we saw it at the new
Cinema in the city center, the actors were awesome.





The movie was great we saw it at the new
Cinema in the city center, right down by the
restaurant where we went for my birthday that
one year, the one where the clown was
inside the cake, the actors were awesome.





The movie was great, the actors were awesome.



Preprocessing

The movie was great, the actors were awesome.

↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓

["the" "movie" "was" "great" "," "the" "actors" "were" "awesome" "."]

Tokenization



The movie was great, the actors were awesome.

[<s> "the" "movie" "was" "great" "," "the" "actors" "were" "awesome" "." <e>]

Tokenization

The movie was great, the actors were awesome.

[<s> "the" "movie" "was" "great" "," "the" "actors" "were" "awesome" "." <e>]

```
In[1]: from nltk.tokenize import word_tokenize

tokens = word_tokenize(s)
tokens = [token.lower() for token in tokens]
tokens = ['<s>'] + tokens + ['<e>']
print(tokens)
```

```
Out[1]: ['<s>', 'the', 'movie', 'was', 'great', ',', 'the', 'actors',
'were', 'awesome', '.', '<e>']
```

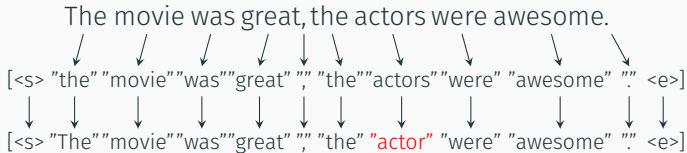
Tokenization

Preprocessing

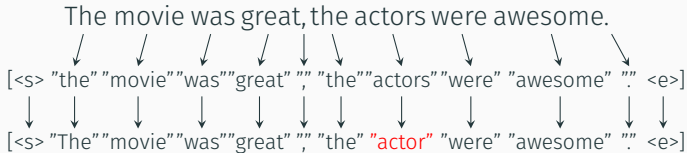
The movie was great, the actors were awesome.

[<s> "the" "movie" "was" "great" "," "the" "actors" "were" "awesome" "." <e>]

[<s> "The" "movie" "was" "great" "," "the" "actor" "were" "awesome" "." <e>]



Stemming



```
In[1]: from nltk.stem.snowball import SnowballStemmer

        stemmer = SnowballStemmer('english')
        stemmed = [stemmer.stem(token) for token in tokens]
        stemmed
```

```
Out[1]: ['<s>', 'the', 'movi', 'was', 'great', ',', 'the', 'actor',
          'were', 'awesom', '.', '<e>']
```

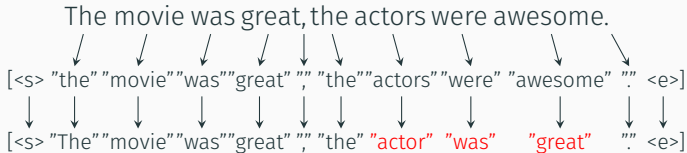
Stemming

The movie was great, the actors were awesome.

[<s> "the" "movie" "was" "great" "," "the" "actors" "were" "awesome" "." <e>]

[<s> "The" "movie" "was" "great" "," "the" "actor" "was" "great" "." <e>]

Lemmatization

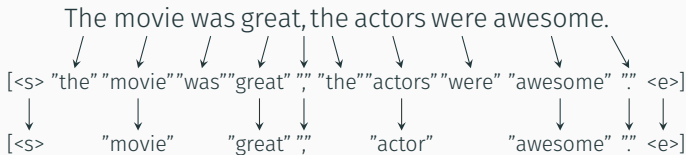


```
In[1]: from nltk.stem import WordNetLemmatizer

lemmatizer = WordNetLemmatizer()
lemmatized = [lemmatizer.lemmatize(token) for token in tokens]
print(lemmatized)
```

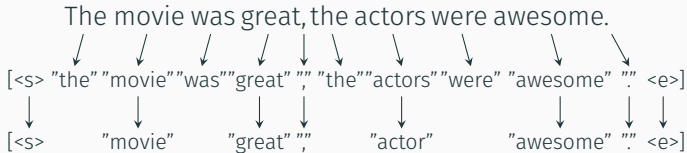
```
Out[1]: ['<s>', 'the', 'movie', 'wa', 'great', ',', 'the', 'actor',
'were', 'awesome', '.', '<e>']
```

Lemmatization



Stopword removal

Preprocessing



```
In[1]: from nltk.corpus import stopwords

pruned = [token for token in tokens if not token in stopwords.
           words('english')]
print(pruned)
```

```
Out[1]: ['<s>', 'movie', 'great', ',', 'actors', 'awesome', '.', '<e>']
```

Stopword removal

Preprocessing

The movie was great, the actors were awesome.

[<s> "the" "movie" "was" "great" "," "the" "actors" "were" "awesome" "." <e>]

["", ".", <e> <s> "actors" "awesome" "great" "movie" "the" "was" "were"]

0 1 2 3 4 5 6 7 8 9 10

Preprocessing

The movie was great, the actors were awesome.

[<s> "the" "movie" "was" "great" "," "the" "actors" "were" "awesome" "." <e>]

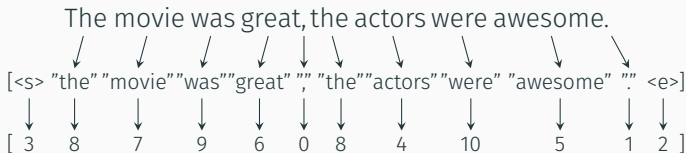
[3 8 7 9 6 0 8 4 10 5 1 2]

[",", ".", <e> <s> "actors" "awesome" "great" "movie" "the" "was" "were"]

0 1 2 3 4 5 6 7 8 9 10

Integer encoding

Preprocessing



Integer encoding

Bag of words



UNIVERSITETET
I OSLO

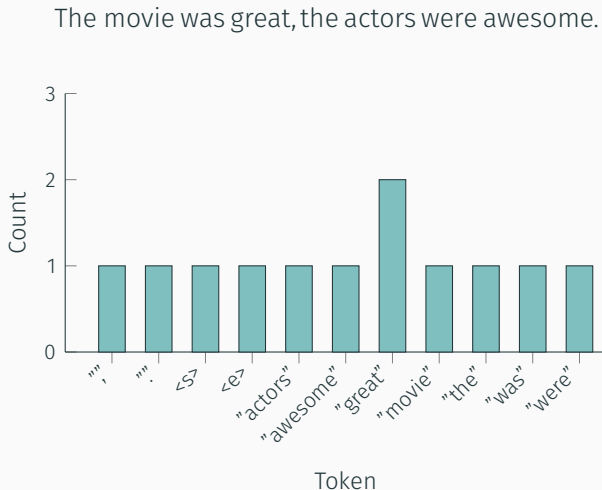
Bag of words



The movie was great, the actors were awesome.



Bag of words





The movie was great, the actors were awesome.

,	.	<s>	<e>	actors	awesome	great	movie	the	was	were
1	1	1	1	1	1	2	1	1	1	1





The movie was great, the actors were awesome.



,	.	<s>	<e>	actors	awesome	awful	great	horrible	movie	the	was	were	sentiment
1	1	1	1	1	1	0	1	0	1	2	1	1	positive
1	1	1	1	1	0	1	0	1	1	2	1	1	negative

The movie was awful, the actors were horrible.





The movie was great, the actors were awesome.

,	.	<s>	<e>	actors	awesome	awful	great	horrible	movie	the	was	were	sentiment
1	1	1	1	1	1	0	1	0	1	2	1	1	positive
1	1	1	1	1	0	1	0	1	1	2	1	1	negative

The movie was awful, the actors were horrible.



Bag of words



	,	.	<s>	<e>	actors	awesome	awful	great	horrible	movie	the	was	were	sentiment
1	1	1	1	1	1	1	0	1	0	1	2	1	1	positive
1	1	1	1	1	1	0	1	0	1	1	2	1	1	negative

$$y = \beta_0 + \sum_i \beta_i X_i$$



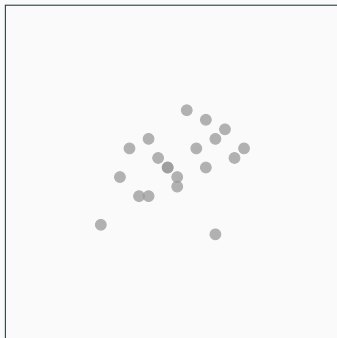
Bag of words



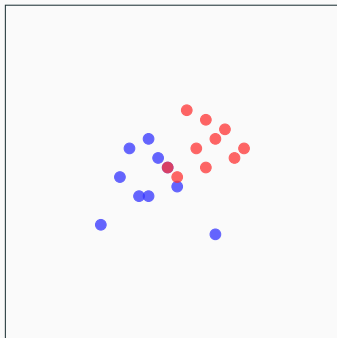
<http://localhost:8888/notebooks/notebooks/Bag%20of%20words%20demo.ipynb>



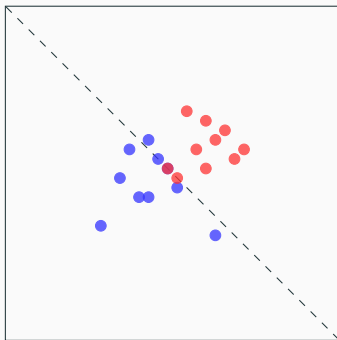
Bag of words: Disadvantages



Bag of words: Disadvantages



Bag of words: Disadvantages



Bag of words: Disadvantages



Dataset: ["This is awesome", "This is wonderful"]



Bag of words: Disadvantages



Dataset: ["This is awesome", "This is wonderful"]

Tokens: [["this" "is" "awesome"], ["this" "is" "wonderful"]]



Bag of words: Disadvantages



Dataset: ["This is awesome", "This is wonderful"]

Tokens: [["this" "is" "awesome"], ["this" "is" "wonderful"]]

Pruned: [["awesome"], ["wonderful"]]



Bag of words: Disadvantages



Dataset: ["This is awesome", "This is wonderful"]

Tokens: [["this" "is" "awesome"], ["this" "is" "wonderful"]]

Pruned: [["awesome"], ["wonderful"]]

Dictionary: ["awesome", "wonderful"]



Bag of words: Disadvantages



Dataset: ["This is awesome", "This is wonderful"]

Tokens: [["this" "is" "awesome"], ["this" "is" "wonderful"]]

Pruned: [["awesome"], ["wonderful"]]

Dictionary: ["awesome", "wonderful"]

Encoded:

	awesome	wonderful
	1	0
	0	1



Bag of words: Disadvantages



Dataset: ["This is awesome", "This is wonderful"]

Tokens: [["this" "is" "awesome"], ["this" "is" "wonderful"]]

Pruned: [["awesome"], ["wonderful"]]

Dictionary: ["awesome", "wonderful"]

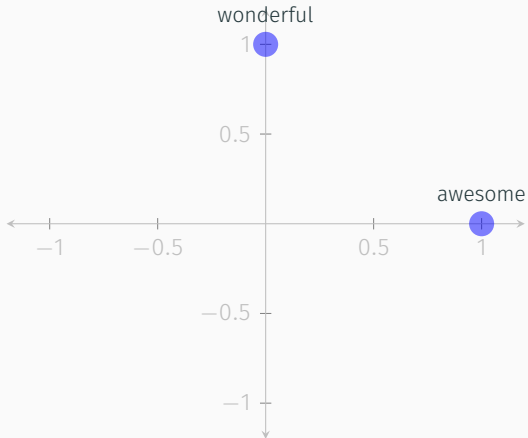
Encoded:

	awesome	wonderful
	1	0
	0	1

Vectors: [[1, 0], [0, 1]]



Bag of words: Disadvantages

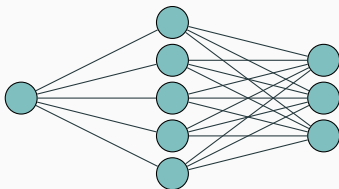


Embeddings

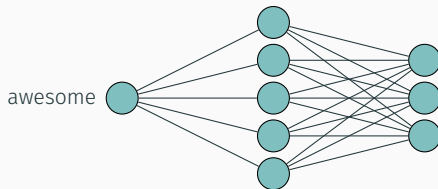


UNIVERSITETET
I OSLO

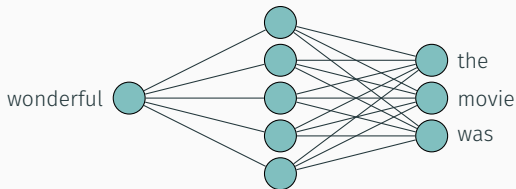
The movie was awesome.
The movie was wonderful.
The movie was fantastic.



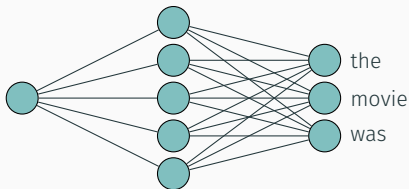
The movie was awesome.
The movie was wonderful.
The movie was fantastic.



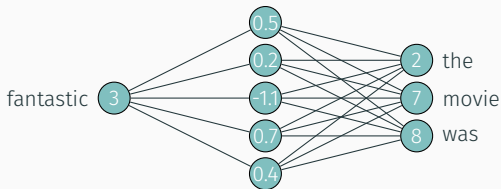
The movie was awesome.
The movie was wonderful.
The movie was fantastic.



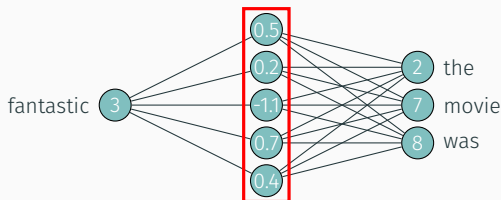
The movie was awesome.
The movie was wonderful.
The movie was fantastic.



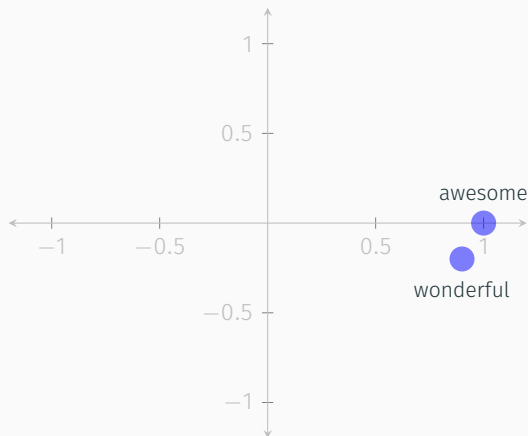
The movie was awesome.
The movie was wonderful.
The movie was fantastic.



The movie was awesome.
The movie was wonderful.
The movie was fantastic.



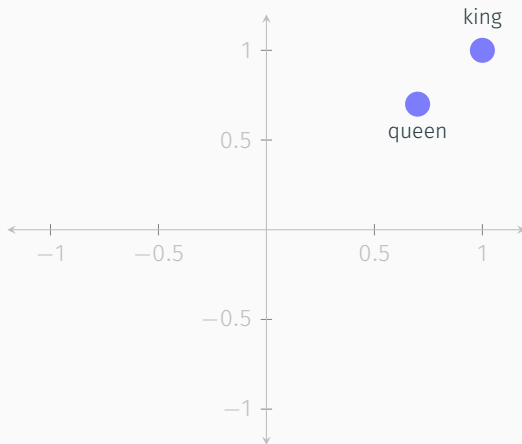
fantastic=[0.5, 0.2, -1.1, 0.7, 0.4]

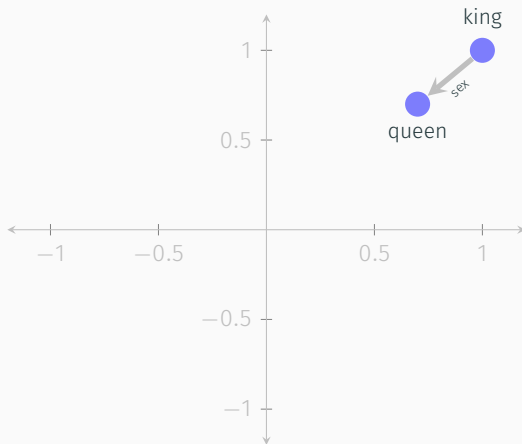


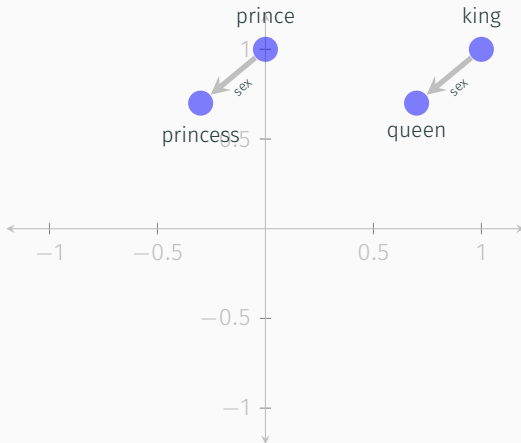


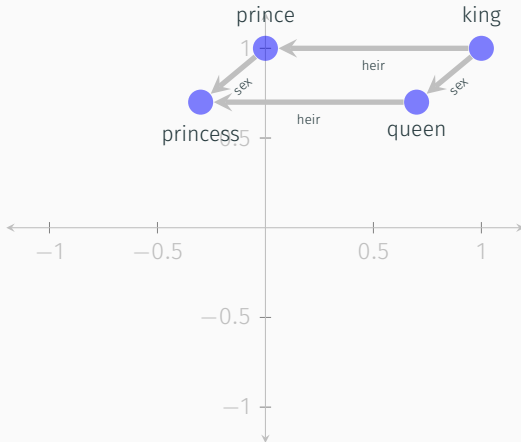
The movie was awesome.
The food was awesome.
The book was awesome.







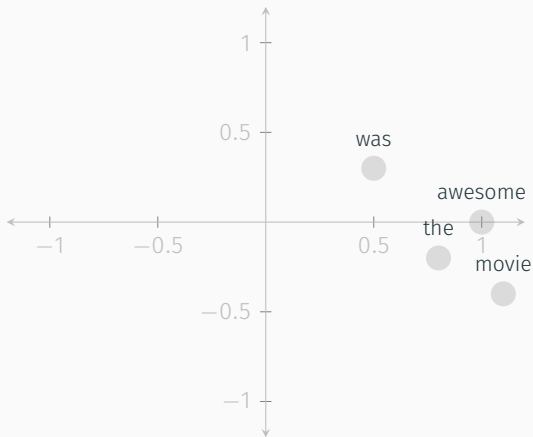


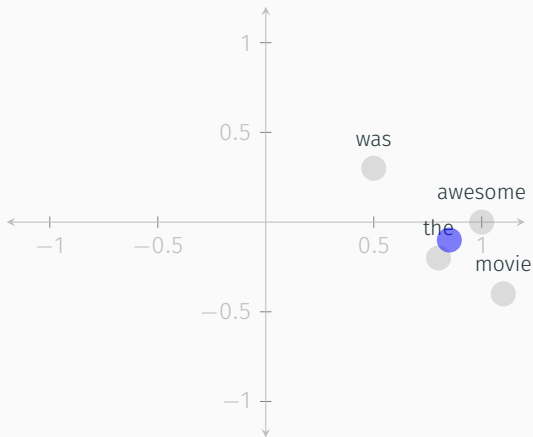


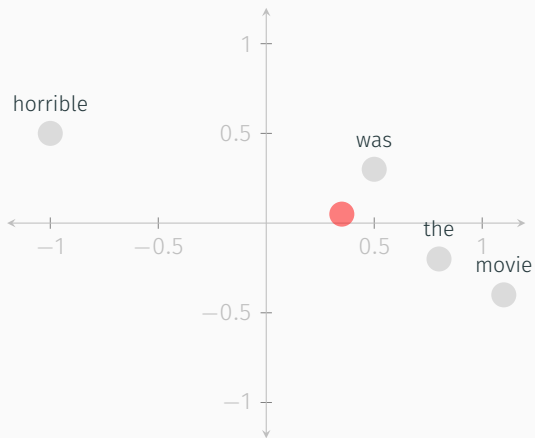


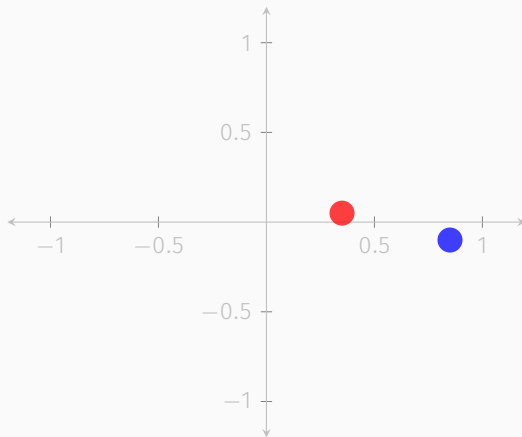
$\text{word2vec}(\text{queen}) = \text{word2vec}(\text{king})$
 $\quad - \text{word2vec}(\text{man})$
 $\quad + \text{word2vec}(\text{woman})$













<http://localhost:8888/notebooks/notebooks/Word2vec%20demo.ipynb>



Word2vec: Disadvantages



Word2vec: Disadvantages

I think the movie was really bad, but my friend said it was good.

=

I think the movie was really good, but my friend said it was bad.

Recurrent neural networks



UNIVERSITETET
I OSLO

Transformers



UNIVERSITETET
I OSLO