



NLP Lab 1

Word Embeddings

# Word Embeddings

## Motivation

- Sprache muss irgendwie in den Computer kommen.
  - Wort/Laut => Zahl(en)
- Das Vokabular einer Sprache ist in der Regel sehr umfangreich.
  - Dimensionalitätsproblem (?)
- Einige Wörter sind enger verbunden als andere
  - Ähnlichkeit, Abstand => Maße?

# Word Embeddings

## Wortvektoren

Katze:	[	1	0	0	]
Hund:	[	0	1	0	]
Maus:	[	0	0	1	]

- One-hot Vektoren
  - hochdimensional, linear unabhängig
- Word Embeddings:
  - reellwertige Vektoren mit  $N \ll |V|$  Dimensionen
  - Komponente im Vektor bildet Eigenschaft des Wortes ab
- wichtige Publikation:
  - Efficient Estimation of Word Representations in Vector Space – Mikolov et al. 2013

# Word Embeddings

## Grundidee

Bellend rennt ... der Katze hinterher.

# Word Embeddings

## Funktionsweise

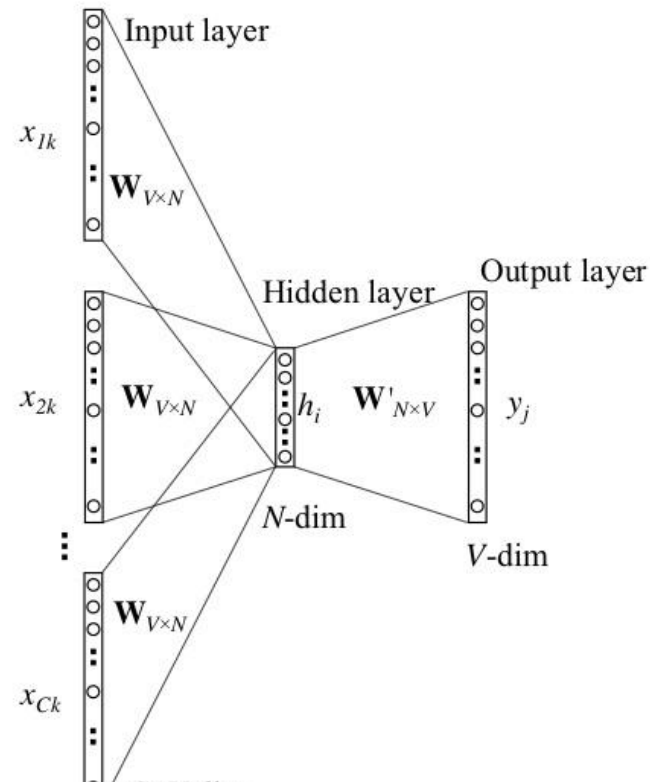
- Wie kann man Wörter auf Vektoren abbilden?
- Trick:
  - Fake Task
  - (einfaches) neuronales Netz, das Wörter aus ihrem Kontext vorhersagt
  - Vektoren werden aus den Gewichtsmatrizen des neuronalen Netz extrahiert
  - Ergebnisse interessieren eigentlich nicht

# Word Embeddings

## Continuous Bag of Words Model (CBOW)

Kontext => fehlendes Wort

- Input Layer:
  - Kontextwörter als Vektoren
- Hidden Layer:
  - N Dimensionen
  - Durchschnitt über Inputvektoren
- Output layer:
  - Wort als Vektor
- Activation: Softmax

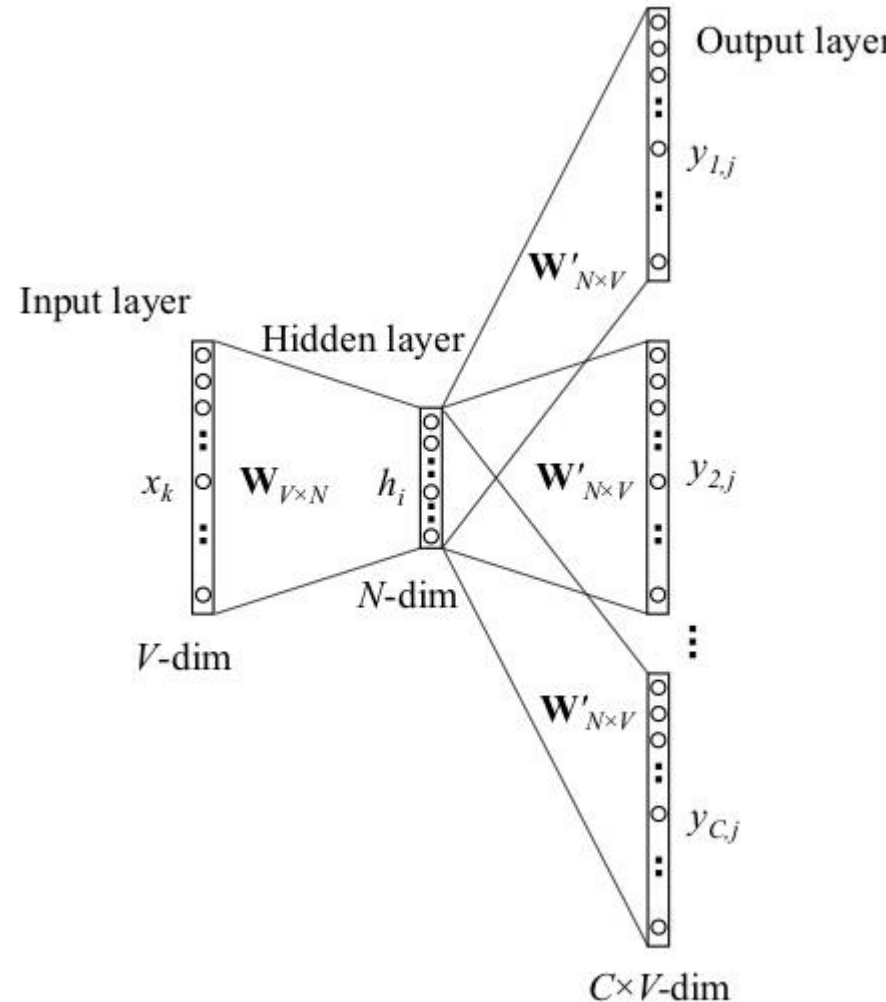


# Word Embeddings

## Skip-Gram Model

Wort => Kontext

- Input
  - Wort als Vektor
- Hidden layer:
  - N Dimensionen
- Output:
  - $C \times V$  dimensionaler Vektor
- Activation: Softmax

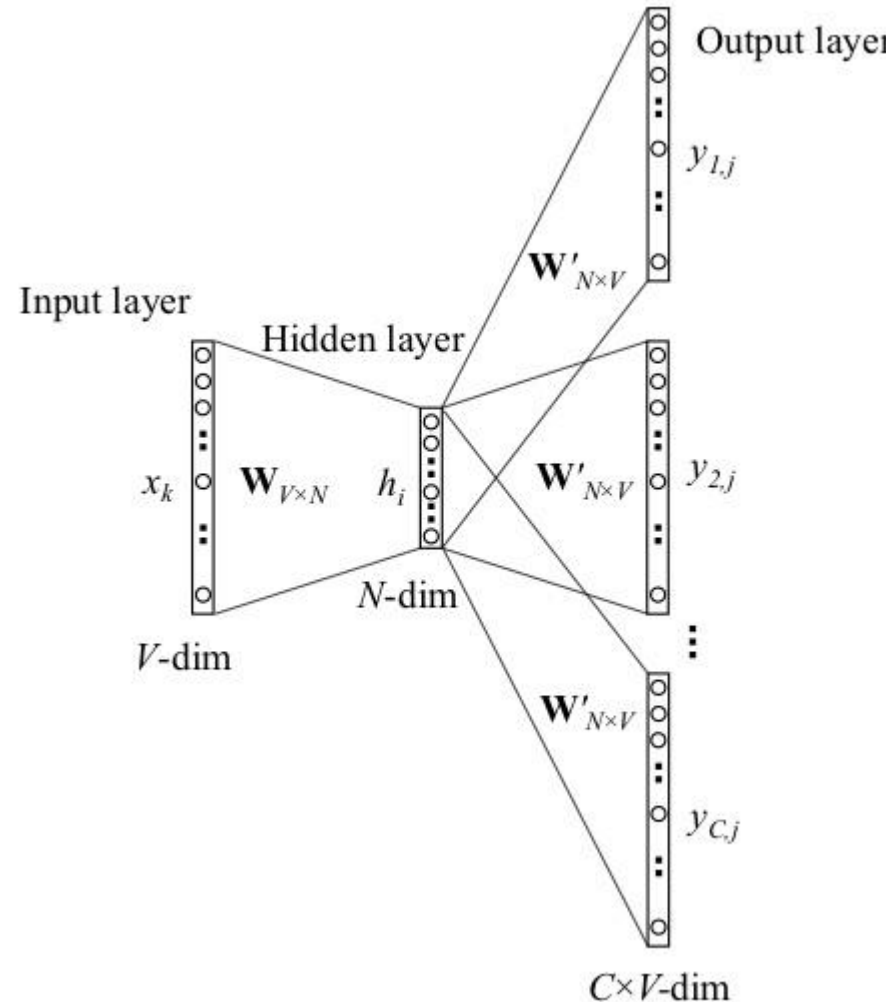


# Word Embeddings

## Skip-Gram Model

Wort => Kontext

- Input
  - Wort als Vektor
- Hidden layer:
  - N Dimensionen
- Output:
  - $C \times V$  dimensionaler Vektor
- Activation: Softmax





# Word Embeddings

## Implementierungen

- GloVe
  - <https://nlp.stanford.edu/projects/glove> (Stanford)
- FastText
  - <https://fasttext.cc> (Facebook)
- Word2Vec
  - <https://code.google.com/archive/p/word2vec> (Google)
- Gensim
  - <https://radimrehurek.com/gensim/models/word2vec.html>

# Word Embeddings — Hands on

# Word Embeddings

## Laborinhalte

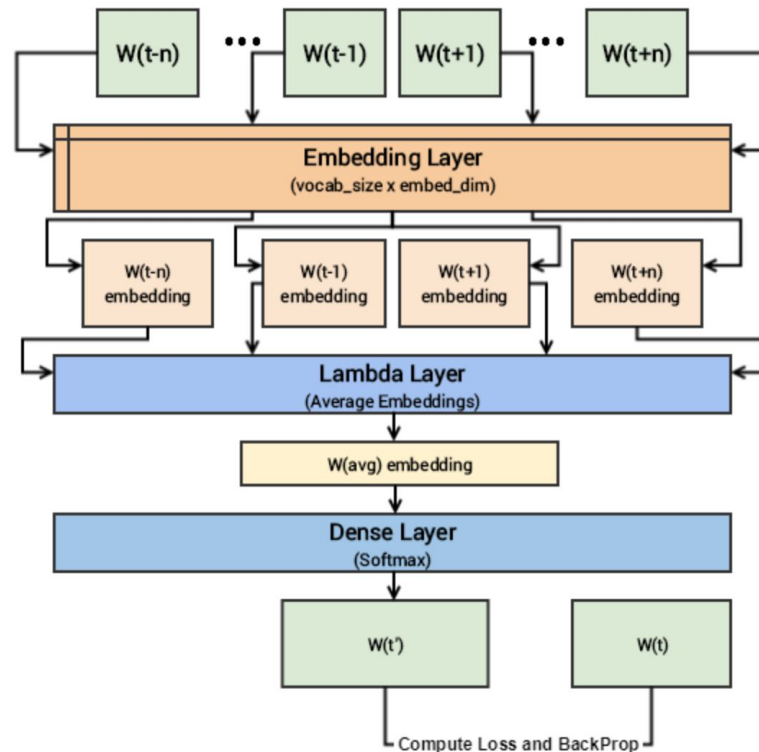
- Implementierung von CBOW in Keras
  - basal, nicht performanceoptimiert
- Arbeit mit Gensim
- Word Embeddings und Bias

# Word Embeddings

## CBOW-Implementierung

Inspiration: [Blogpost von Dipanjan Sarkar](#)

- [Embedding Layer](#) dient als "Lookup-Table"
- Lambda-Layer, um Kontextvektoren zu mitteln
- Dense Layer erzeugt per Softmax-Aktivierung Vorhersagen des fehlenden Wortes



Visual depiction of the CBOW deep learning model

# Word Embeddings

## Bias

- [What are the biases in my word embedding?](#)
- [Understanding the Origins of Bias in Word Embeddings](#)
- [Man is to Computer Programmer as Woman is to Homemaker? Debiasing Word Embeddings](#)
- [Evaluating the Underlying Gender Bias in Contextualized Word Embeddings](#)

<https://github.com/hskaailabnlp/embeddings>