

# Word2Vec und das Skip-Gram-Modell von Grund auf skripten: Eigene Word Embeddings erstellen

Aleksandr Schamberger

Humboldt-Universität zu Berlin

Institut für Slawistik und Hungarologie

Sprachenübergreifend: Computerlinguistik II – Digitale Sprachmodelle und ihre  
Anwendung

Sommersemester 2024

2. Juli 2024

# Gliederung

- 1 Wörter (sinnvoll) als Zahlen darstellen
- 2 Worteinbettungen (Word Embeddings)

# Gliederung

1 Wörter (sinnvoll) als Zahlen darstellen

2 Worteinbettungen (Word Embeddings)

# Wörter (sinnvoll) als Zahlen darstellen

- Bank
- Stuhl
- Sitzgelegenheit
- dick
- dünn
- schlank

# ASCII-Tabelle

Scan- code	ASCII hex dez	Zeichen	Scan- code	ASCII hex dez	Zch.	Scan- code	ASCII hex dez	Zch.	Scan- code	ASCII hex dez	Zch.
	00 0	NUL ^@		20 32	SP		40 64	@	0D	60 96	`
	01 1	SOH ^A	02	21 33	!	1E	41 65	A	1E	61 97	a
	02 2	STX ^B	03	22 34	"	30	42 66	B	30	62 98	b
	03 3	ETX ^C	29	23 35	#	2E	43 67	C	2E	63 99	c
	04 4	EOT ^D	05	24 36	\$	20	44 68	D	20	64 100	d
	05 5	ENQ ^E	06	25 37	%	12	45 69	E	12	65 101	e
	06 6	ACK ^F	07	26 38	&	21	46 70	F	21	66 102	f
	07 7	BEL ^G	0D	27 39	'	22	47 71	G	22	67 103	g
0E	08 8	BS ^H	09	28 40	(	23	48 72	H	23	68 104	h
0F	09 9	TAB ^I	0A	29 41	)	17	49 73	I	17	69 105	i
	0A 10	LF ^J	1B	2A 42	*	24	4A 74	J	24	6A 106	j
	0B 11	VT ^K	1B	2B 43	+	25	4B 75	K	25	6B 107	k
	0C 12	FF ^L	33	2C 44	,	26	4C 76	L	26	6C 108	l
1C	0D 13	CR ^M	35	2D 45	-	32	4D 77	M	32	6D 109	m
	0E 14	SO ^N	34	2E 46	.	31	4E 78	N	31	6E 110	n
	0F 15	SI ^O	08	2F 47	/	18	4F 79	O	18	6F 111	o
	10 16	DLE ^P	0B	30 48	0	19	50 80	P	19	70 112	p
	11 17	DC1 ^Q	02	31 49	1	10	51 81	Q	10	71 113	q
	12 18	DC2 ^R	03	32 50	2	13	52 82	R	13	72 114	r
	13 19	DC3 ^S	04	33 51	3	1F	53 83	S	1F	73 115	s
	14 20	DC4 ^T	05	34 52	4	14	54 84	T	14	74 116	t
	15 21	NAK ^U	06	35 53	5	16	55 85	U	16	75 117	u
	16 22	SYN ^V	07	36 54	6	2F	56 86	V	2F	76 118	v
	17 23	ETB ^W	08	37 55	7	11	57 87	W	11	77 119	w
	18 24	CAN ^X	09	38 56	8	2D	58 88	X	2D	78 120	x
	19 25	EM ^Y	0A	39 57	9	2C	59 89	Y	2C	79 121	y
	1A 26	SUB ^Z	34	3A 58	:	15	5A 90	Z	15	7A 122	z
01	1B 27	Esc ^[	33	3B 59	;		5B 91	[		7B 123	{
	1C 28	FS ^\	2B	3C 60	<		5C 92	\		7C 124	
	1D 29	GS ^]	0B	3D 61	=		5D 93	]		7D 125	}
	1E 30	RS ^^	2B	3E 62	>	29	5E 94	^		7E 126	~
	1F 31	US ^_	0C	3F 63	?	35	5F 95	_	53	7F 127	DEL

# Wörter als deren ASCII-Zeichenkodierung

- Bank
  - 66 97 110 107 10
- Stuhl
  - 83 116 117 104 108
- Sitzgelegenheit
  - 83 105 116 122 103 101 108 101 103 101 110 104 101 105 116
- dick
  - 100 105 99 107
- dünn
  - 100 195 188 110 110
- schlank
  - 115 99 104 108 97 110 107

# Semantische Relationen (und Eigenschaften)

- **Hyponymie/Hyperonymie:** (Bank,Stuhl) → Sitzgelegenheit
- **Ko-Hyponymie:** Bank - Stuhl
- **Antonymie:** dünn - dick
- **Synonymie:** dünn - schlank
- **Ambiguität:** Bank - Bank

# Semantische Relationen (und Eigenschaften)

- **Hyponymie/Hyperonymie:** (Bank,Stuhl) → Sitzgelegenheit
- **Ko-Hyponymie:** Bank - Stuhl
- **Antonymie:** dünn - dick
- **Synonymie:** dünn - schlank
- **Ambiguität:** Bank - Bank

Und viele weitere, sprachliche Eigenschaften ...



# Zwischenlösung: One-Hot-Vektoren

- Bank
  - 1 0 0 0 0 0
- Stuhl
  - 0 1 0 0 0 0
- Sitzgelegenheit
  - 0 0 1 0 0 0
- dick
  - 0 0 0 1 0 0
- dünn
  - 0 0 0 0 1 0
- schlank
  - 0 0 0 0 0 1

# Problem: Der Umfang des deutschen Wortschatzes

Umfang deutscher Wörter (Grundformen)

# Problem: Der Umfang des deutschen Wortschatzes

Umfang deutscher Wörter (Grundformen)

- Aktiver Wortschatz: 12.000 - 16.000

# Problem: Der Umfang des deutschen Wortschatzes

Umfang deutscher Wörter (Grundformen)

- Aktiver Wortschatz: 12.000 - 16.000
- Passiver Wortschatz: 50.000

# Problem: Der Umfang des deutschen Wortschatzes

## Umfang deutscher Wörter (Grundformen)

- Aktiver Wortschatz: 12.000 - 16.000
- Passiver Wortschatz: 50.000
- Im Allgemeinen angenommen: 300.000 - 500.000

# Problem: Der Umfang des deutschen Wortschatzes

## Umfang deutscher Wörter (Grundformen)

- Aktiver Wortschatz: 12.000 - 16.000
- Passiver Wortschatz: 50.000
- Im Allgemeinen angenommen: 300.000 - 500.000
- Dudenkorpus: > 18 Millionen

# Problem: Der Umfang des deutschen Wortschatzes

## Umfang deutscher Wörter (Grundformen)

- Aktiver Wortschatz: 12.000 - 16.000
- Passiver Wortschatz: 50.000
- Im Allgemeinen angenommen: 300.000 - 500.000
- Dudenkorpus: > 18 Millionen

10.000 Wörter: Bank  $\rightarrow 1_1 0_2 0_3 0_4 0_5 0_6 0_7 0_8 0_9 \dots 0_{10.000}$

# Gliederung

- 1 Wörter (sinnvoll) als Zahlen darstellen
- 2 Worteinbettungen (Word Embeddings)



# Von weiten Vektoren zu dichten Vektoren

10.000 Wörter: Bank  $\rightarrow 1_1 0_2 0_3 0_4 0_5 0_6 0_7 0_8 0_9 \dots 0_{10.000}$

# Von weiten Vektoren zu dichten Vektoren

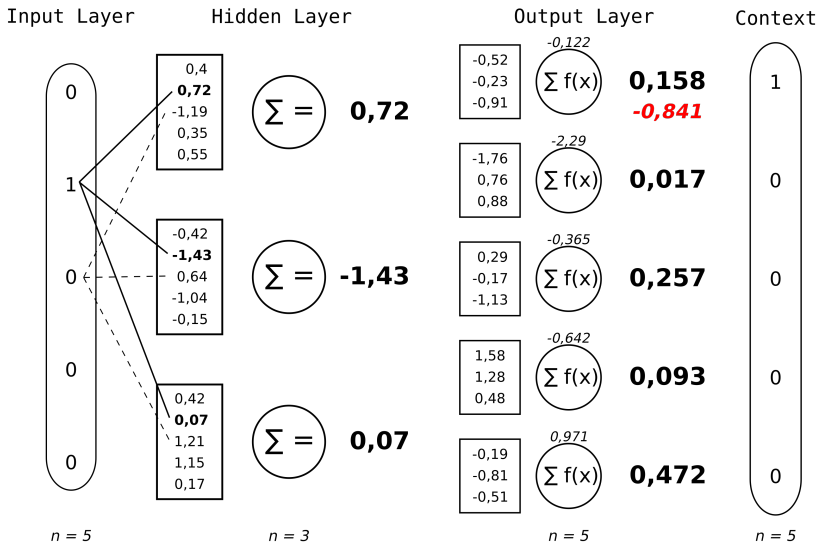
10.000 Wörter: Bank  $\rightarrow 1_1 0_2 0_3 0_4 0_5 0_6 0_7 0_8 0_9 \dots 0_{10.000}$

$\leftrightarrow$

10.000 Wörter: Bank  $\rightarrow$

$-1, 23_1 3, 4_2 0, 09_3 0_4 1, 15_5 0, 98_6 - 2, 34_7 0, 11_8 - 1, 3_9 \dots 0_{300}$

# Skip-Gram-Modell: Gesamtmodell mit Beispielrechnung



# Skip-Gram-Modell: Kontext (Syntagma)

## Source Text

## Training Samples

The quick brown fox jumps over the lazy dog. →

(the, quick)  
(the, brown)

The quick brown fox jumps over the lazy dog. →

(quick, the)  
(quick, brown)  
(quick, fox)

The quick brown fox jumps over the lazy dog. →

(brown, the)  
(brown, quick)  
(brown, fox)  
(brown, jumps)

The quick brown fox jumps over the lazy dog. →

(fox, quick)  
(fox, brown)  
(fox, jumps)  
(fox, over)

# Quellen

- Tae, Jake (13.07.2020): Word2vec from Scratch  
<<https://jaketae.github.io/study/word2vec/>> (abgerufen am 02.07.2024)
- McCormick, Chris (19.04.2016): Word2Vec Tutorial - The Skip-Gram Model  
<<http://mccormickml.com/2016/04/19/word2vec-tutorial-the-skip-gram-model/>> (abgerufen am 02.07.2024)
- Tae, Jake (05.02.2020): Building Neural Network From Scratch <<https://jaketae.github.io/study/neural-net/>> (abgerufen am 02.07.2024)
- Tae, Jake (21.12.2019): Demystifying Entropy (And More)  
<<https://jaketae.github.io/study/information-entropy/>> (abgerufen am 02.07.2024)

**Beispiel(e) im Editor!**