

Sentimentanalyse Harry Potter

Projektarbeit Deep Learning

Emilia Annacker

30. Mai 2025

Motivation

Beschreibung der Daten

LSTM-Modell mit Twitterdatensatz

BERT mit GoEmotions-Datensatz

Analyse einzelner Sätze

Analyse mit Kontext

Fazit

Motivation

Beschreibung der Daten

LSTM-Modell mit Twitterdatensatz

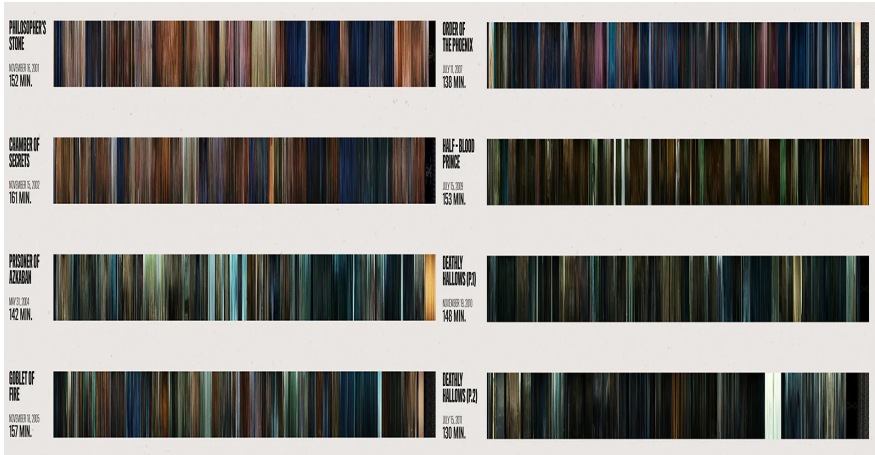
BERT mit GoEmotions-Datensatz

Analyse einzelner Sätze

Analyse mit Kontext

Fazit

Movie Bar Code Plots der Harry-Potter-Filme



Hypothese

- ▶ Sentimentanalyse der Bücher sollte ebenfalls immer düsterere Stimmung im Verlauf der Reihe zeigen
- ▶ Hauptpersonen werden älter
- ▶ zunehmende Bedrohung durch den Antagonisten (vor allem ab Teil 5)



Motivation

Beschreibung der Daten

LSTM-Modell mit Twitterdatensatz

BERT mit GoEmotions-Datensatz

Analyse einzelner Sätze

Analyse mit Kontext

Fazit

Harry-Potter-Bände von J.K. Rowling

Jahr	Titel	Anzahl Sätze
1997	Harry Potter and the Philosopher's Stone	6422
1998	Harry Potter and the Chamber of Secrets	6658
1999	Harry Potter and the Prisoner of Azkaban	8663
2000	Harry Potter and the Goblet of Fire	14306
2003	Harry Potter and the Order of the Phoenix	14112
2005	Harry Potter and the Half-Blood Prince	11919
2007	Harry Potter and the Deathly Hallows	14506
		76586

<https://www.kaggle.com/datasets/moxxis/harry-potter-lstm>

Datenaufbereitung

- ▶ Kapitel unterteilen, durch Kapitelüberschriften den Büchern zuordnen
- ▶ Dataframes buchweise abspeichern
- ▶ Spalten: chapter, sentence

```
      chapter      sentence
0  THE BOY WHO LIVED  Mr and Mrs Dursley , of number four , Privet D...
1  THE BOY WHO LIVED  They were the last people you'd expect to be i...
2  THE BOY WHO LIVED  Mr Dursley was the director of a firm called G...
3  THE BOY WHO LIVED  He was a big , beefy man with hardly any neck ...
4  THE BOY WHO LIVED  Mrs Dursley was thin and blonde and had nearly...
      chapter      sentence
14501 NINETEEN YEARS LATER      " He'll be all right , " murmured Ginny .
14502 NINETEEN YEARS LATER  As Harry looked at her , he lowered his hand a...
14503 NINETEEN YEARS LATER      " I know he will . "
14504 NINETEEN YEARS LATER  The scar had not pained Harry for nineteen yea...
14505 NINETEEN YEARS LATER      All was well .
```


Motivation

Beschreibung der Daten

LSTM-Modell mit Twitterdatensatz

BERT mit GoEmotions-Datensatz

Analyse einzelner Sätze

Analyse mit Kontext

Fazit

Training mit Twitterdatensatz

- ▶ 1048573 Tweets
- ▶ Labels: positiv/negativ
- ▶ TextVectorization: Punctuation entfernen, lower, split='whitespace'
- ▶ Modell mit Embeddingschicht (3 Dimensionen), 2 LSTM-Schichten und 2 Dense-Schichten

Model: "sequential_2"

Layer (type)	Output Shape	Param #
embedding (Embedding)	(None, 25, 3)	0
lstm (LSTM)	(None, 25, 100)	41,500
lstm_1 (LSTM)	(None, 100)	80,400
dense (Dense)	(None, 100)	10,100
dense_1 (Dense)	(None, 1)	100

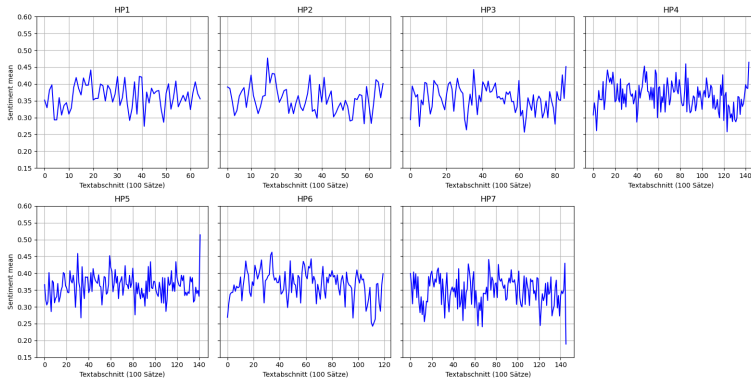
Total params: 132,000 (516.43 KB)
Trainable params: 132,000 (516.43 KB)
Non-trainable params: 0 (0.00 B)

```
model.fit(X_train_tok, y_train, epochs=5, batch_size=100, verbose=True,  
validation_data=(X_test_tok, y_test))
```

```
model.compile(optimizer=keras.optimizers.Adam(learning_rate=0.01),  
loss='binary_crossentropy',  
metrics=['accuracy'])
```

Ergebnisse

- ▶ min: 0.000000094
- ▶ max: 0.998
- ▶ mean: 0.359



Motivation

Beschreibung der Daten

LSTM-Modell mit Twitterdatensatz

BERT mit GoEmotions-Datensatz

Analyse einzelner Sätze

Analyse mit Kontext

Fazit

BERT

- ▶ Language Model, Google, 2018
 - ▶ bidirektionaler Transformer, Encoder-only-Transformer-Architektur
 - ▶ vortrainiert auf ungelabeltem Text: Voraussage des nächsten Tokens durch Masked Tokens und Next-Sentence-Prediction
 - ▶ kontextuelle Wortrepräsentationen
 - ▶ können durch Fine-Tuning eines zusätzlichen Layers oder Heads für bestimmte NLP-Aufgaben angepasst werden kann
- hier Sentimentanalyse

GoEmotions von Google

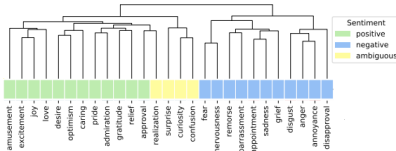
- ▶ manuell annotiertes Dataset zu Emotionen (2020)
- ▶ 58k englisch Reddit-Kommentare
- ▶ Labels: 27 Emotionen und neutral

Number of examples	58,009
Number of emotions	27 + neutral
Number of unique raters	82
Number of raters / example	3 or 5
Marked unclear or difficult to label	1.6%
Number of labels per example	1: 83% 2: 15% 3: 2% 4+: .2%
Number of examples w/ 2+ raters agreeing on at least 1 label	54,263 (94%)
Number of examples w/ 3+ raters agreeing on at least 1 label	17,763 (31%)

Table 2: Summary statistics of our labeled data.

	Text	admiration	amusement	anger	annoyance	approval	caring	confusion	curiosity	desire	love	nervousness	optimism	pride	realization	relief	remorse	sadness	surprise	neutral
0	My favourite food is anything I didn't have to...	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1
1	Now if he does off himself, everyone will thin...	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1
2	WHY THE FUCK IS BAYLESS ISOING	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
3	To make her feel threatened	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
4	Dirty Southern Wankers	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

GoEmotions



- ▶ die 27 Emotionen können hierarchisch geclustert werden
- ▶ ambige Emotionen häufiger in positiven Kontexten

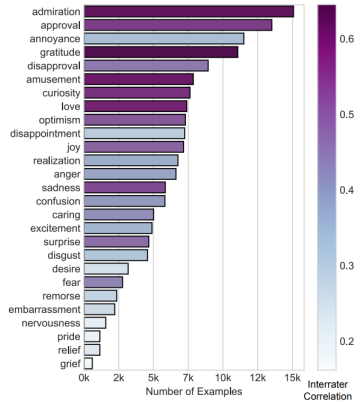


Figure 1: Our emotion categories, ordered by the number of examples where at least one rater uses a particular label. The color indicates the interrater correlation.

Fine-Tuning des Modells

Modell: <https://huggingface.co/bhadresh-savani/bert-base-go-emotion>

- ▶ Basismodell: distilbert-base-uncased, kleinere schnellere Version von BERT
- ▶ Klassifikationshead: Multilabel-Sequence-Classification von 28 Emotionen
- ▶ 6 Transformer-Layer
- ▶ epochs = 3
- ▶ batch size = 16

```
tokenizer =  
AutoTokenizer.from_pretrained('bhadresh-savani/bert-base-go-emotion')  
  
model =  
TFBertForSequenceClassification.from_pretrained("bhadresh-savani/  
bert-base-go-emotion", from_pt=True)
```

- ▶ normalisierter Output der 28 Kategorien in Dictionarys

Beispiele

'A horrible , half sucking , half moaning sound came out of the square hole , along with an unpleasant smell like open drains'

{*'annoyance'*: 0.157, *'disgust'*: 0.364, *'fear'*: 0.149}

'And if the Dursleys were unhappy to have him back for the holidays , it was nothing to how Harry felt'

{*'disappointment'*: 0.205, *'sadness'*: 0.282, *'neutral'*: 0.213}

'A braver man than Vernon Dursley would have quailed under the furious look Hagrid now gave him when Hagrid spoke , his every syllable trembled with rage'

{*'amusement'*: 0.402, *'neutral'*: 0.352}

- ▶ Klassifikation satzweise
- ▶ Mittelwert über Kapitel
- ▶ Sortieren bzw. Zusammenfassen der Emotionen

Motivation

Beschreibung der Daten

LSTM-Modell mit Twitterdatensatz

BERT mit GoEmotions-Datensatz

Analyse einzelner Sätze

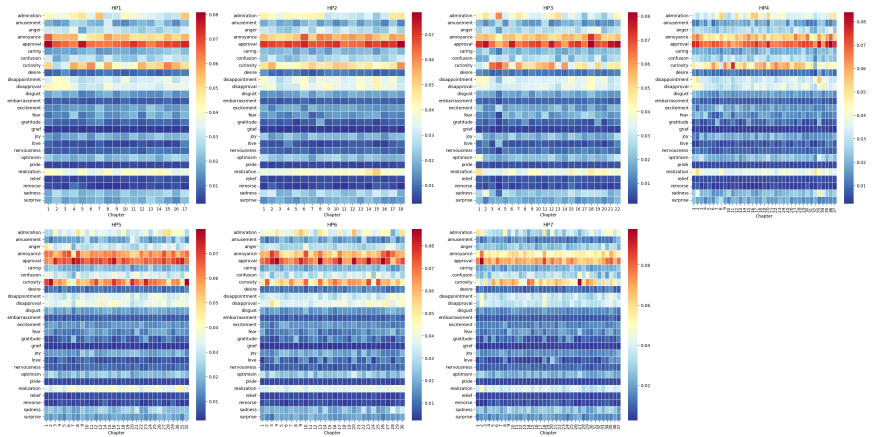
Analyse mit Kontext

Fazit

Ergebnisse I



Ergebnisse II



Zwischenfazit

- ▶ keine erkennbare Tendenz, alle Emotionen buchübergreifend recht gleichmäßig verteilt
- ▶ Problem: jeder Satz wird separat klassifiziert
- ▶ für mehr Kontext zur Klassifikation größere Abschnitte
 - ▶ 10 Sätze auf einmal
 - ▶ Mittelwert über Kapitel zur Vergleichbarkeit

Motivation

Beschreibung der Daten

LSTM-Modell mit Twitterdatensatz

BERT mit GoEmotions-Datensatz

Analyse einzelner Sätze

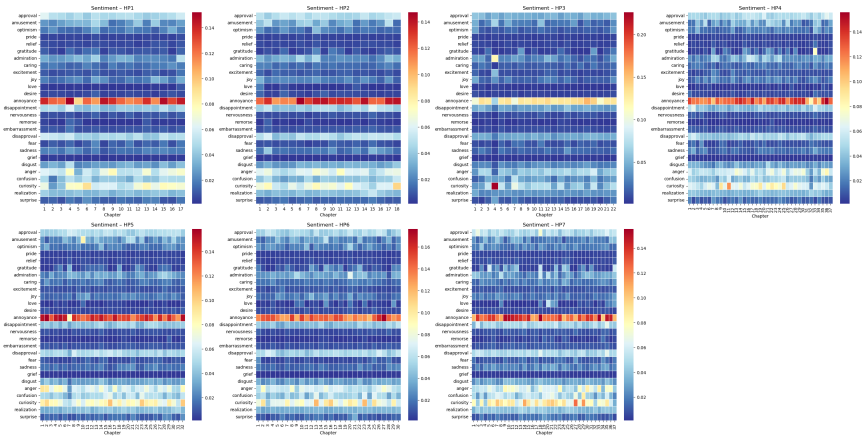
Analyse mit Kontext

Fazit

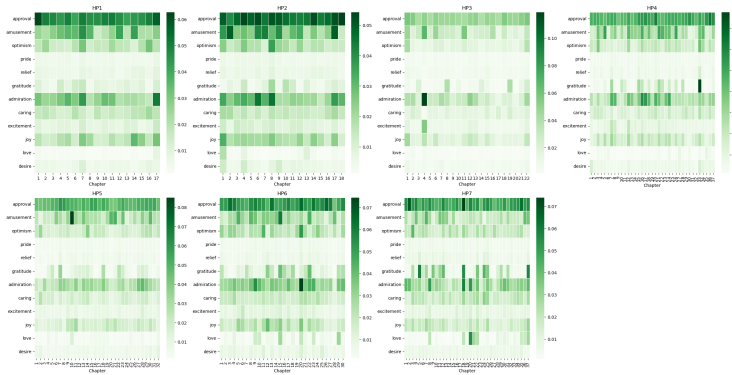
Ergebnisse I



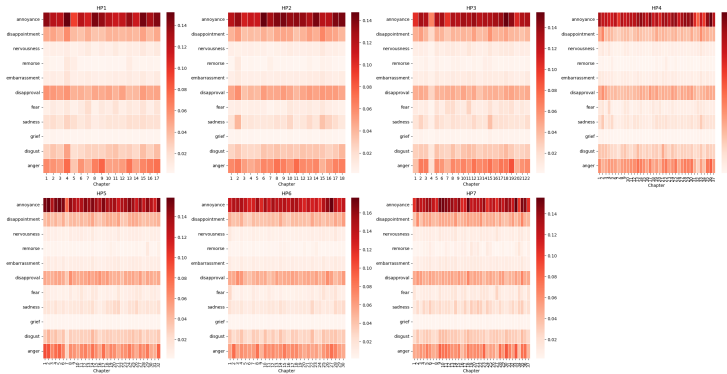
Ergebnisse II



Ergebnisse III

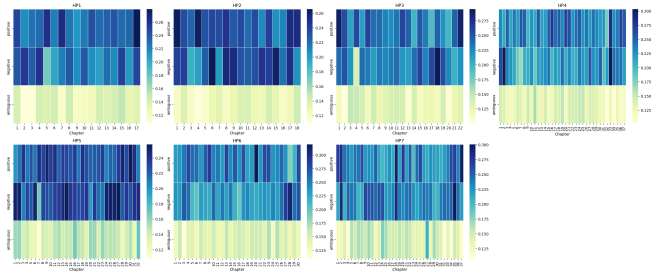


Ergebnisse IV

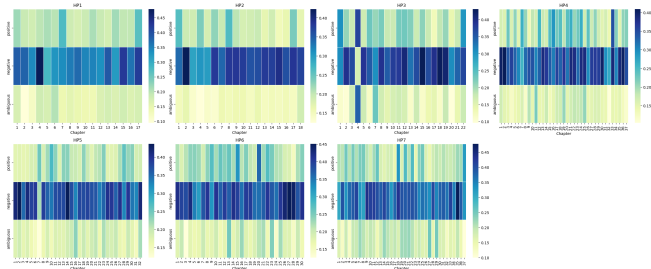


Ergebnisse V

Emotionen satzweise

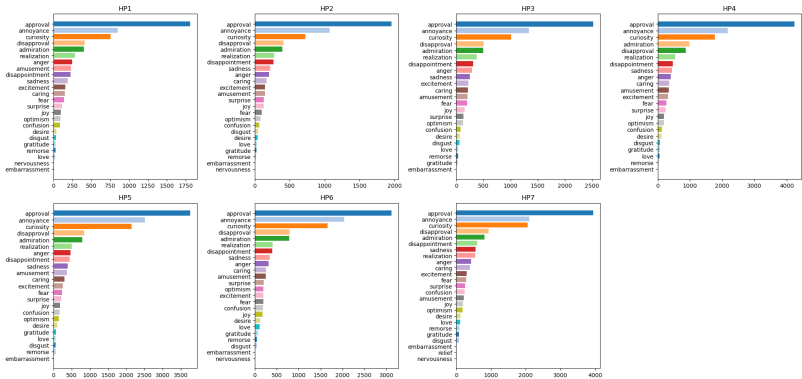


Emotionen im Kontext



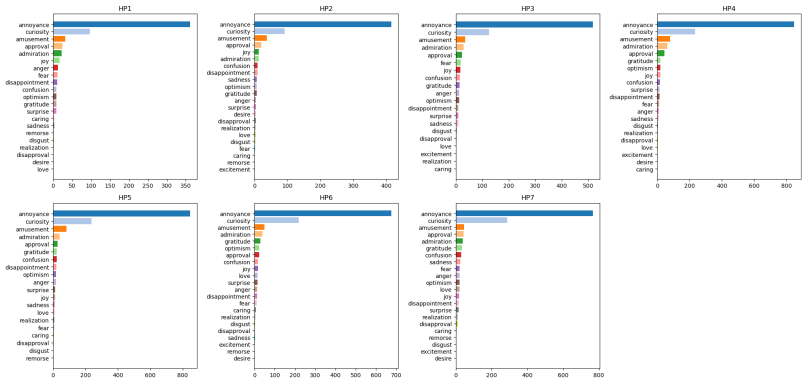
Ergebnisse VI

Häufigkeit der Emotionen satzweise



Ergebnisse VII

Häufigkeit der Emotionen mit Kontext



Motivation

Beschreibung der Daten

LSTM-Modell mit Twitterdatensatz

BERT mit GoEmotions-Datensatz

Analyse einzelner Sätze

Analyse mit Kontext

Fazit

Zusammenfassung

- ▶ bei satzweiser Analyse fast identische Verteilung der Emotionen wie im GoEmotions-Dataset
- ▶ stärkeres Einbeziehen des Kontexts verändert die Klassifikation maßgeblich
- ▶ Training mit HP-Daten könnte hilfreich sein
- ▶ Anwendung von Reddit-Daten auf Literatur:
 - ▶ Reddit: bei GoEmotions wurde die Haltung des Verfassers eingeschätzt
 - ▶ HP: Erzähler beschreibt Emotionen von vielen Personen gleichzeitig, Kontext und Beziehungen entwickeln sich über weite Strecken, Differenzierung der Emotionen funktioniert gar nicht

