

Erkennung von Hate Speech mit Twitter - Zwischenpräsentation

Elena Marion Friedrich, Nasiba Tuychieva, Sven Ole Wall, Imran Nteli Chalil,
Christian Engels

Artificial Intelligence Group,
University of Hagen, Germany

17. Dezember 2024



- 1 Fragestellung
- 2 Workflow
- 3 Close-up: Datenbereinigung
- 4 Close-up: Vektorisierung
- 5 Close-up: Datenexploration
- 6 Besonderheiten des Datensatzes
- 7 Verfeinerung Fragestellung
- 8 Ausblick: Model Training

Ausgangsfrage

Wie kann mithilfe von klassischen Verfahren des maschinellen Lernens sowie mit Deep-Learning-Verfahren Hate Speech auf sozialen Plattformen anhand des Beispiels Twitter erkannt werden?

Klassifikation von Text / Sentiment Analysis

Phase I : Zentralisiert

Datenverständnis und
Datenbereinigung



Datenexploration



Verfeinerung
Forschungsfrage



Vektorisierung



Phase II : Parallelisiert

Model Training



Model Evaluation

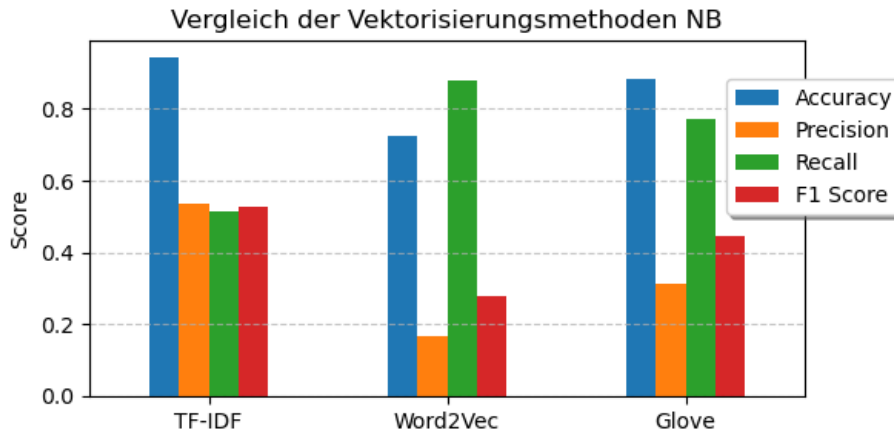
Parameter-Tuning und
Daten-Resampling



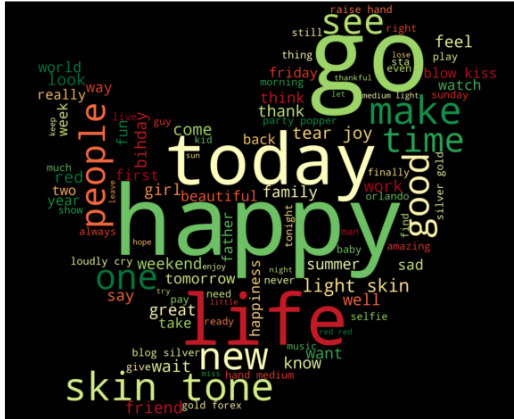
- Entfernung unnötiger Zeichen
- Entfernung von Erwähnungen
- Auftrennen zusammengeschiedener Wörter
- Entfernung von Stopwords und häufigsten Wörtern
- Konvertierung von Emojis zu Text
- Stemming und Lemmatisierung

1. #nationalbestfriendsday #sisters
#family #love thank you @user for
being my best friend ❤️ 😊
2. @user i will be there !!! hoping for a
beautiful day ~ clear skies with a
breeze would be perfect !!! ;)

1. national good friend day sister family
love thank friend sparkle heart
grinning face
2. hope beautiful day ~ clear sky breeze
would perfect



No Hate Speech Wörter bereinigte Daten



Hate Speech Wörter bereinigte Daten

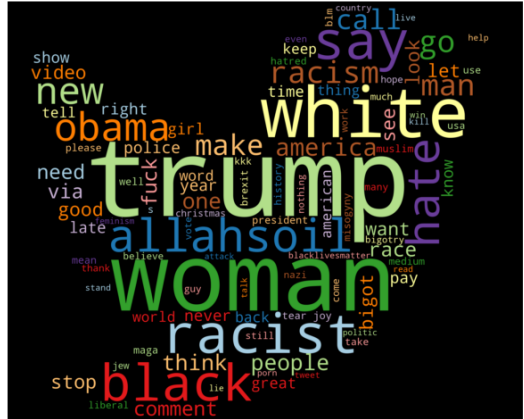


Abbildung: No Hate Speech und Hate Speech Tweets

Binäre Klassifizierung (nur Trainingsdaten)

Labels unterscheiden
zwischen Hassrede und
Nicht-Hassrede

→ Auswahl der Algorithmen und Metriken

Binäre Klassifizierung (nur Trainingsdaten)

Labels unterscheiden
zwischen Hassrede und
Nicht-Hassrede

→ Auswahl der Algorithmen und Metriken

Unausgeglichener Datensatz

- 29720 Datensätze Nicht-Hassrede
- 2242 Datensätze Hassrede

→ Verzerrte Modelle und Metriken

Binäre Klassifizierung (nur Trainingsdaten)

Labels unterscheiden zwischen Hassrede und Nicht-Hassrede

→ Auswahl der Algorithmen und Metriken

Unausgeglichener Datensatz

- 29720 Datensätze Nicht-Hassrede
- 2242 Datensätze Hassrede

→ Verzerrte Modelle und Metriken

Verwendete Sprache

- Rechtschreibfehler
- Abkürzungen
- Emoticons
- @-Erwähnungen
- Hashtags

→ Beeinflussung der Modellgenauigkeit
⇒ Vorverarbeitung

Ausgangsfrage

Wie kann mithilfe von klassischen Verfahren des maschinellen Lernens sowie mit Deep-Learning-Verfahren Hate Speech auf sozialen Plattformen anhand des Beispiels Twitter erkannt werden?

Klassifikation von Text / Sentiment Analysis

Methodische Detailfragen

Welche Möglichkeiten zum Umgang mit Klassenungleichgewichten gibt es und wie ist deren Einfluss auf die Modellperformance?

- Naïve Bayes
- Support-Vektor-Maschine
- Ensemble Learning
- RNN-LSTM
- RNN-GRU



Erkennung von Hate Speech mit Twitter - Zwischenpräsentation

Elena Marion Friedrich, Nasiba Tuychieva, Sven Ole Wall, Imran Nteli Chalil,
Christian Engels

Artificial Intelligence Group,
University of Hagen, Germany

17. Dezember 2024