

Anweisungen für die Markierung von Daten

Hintergrund:

Wir arbeiten daran, Modelle zur automatischen Erkennung von Hass im Internet besser testen zu können. Dafür haben wir einen Datensatz mit ein paar Tausend kurzen Textstücken erstellt. Jeder Eintrag in diesem Datensatz spiegelt einen bestimmten Aspekt von Hassrede wieder.

Deine Aufgaben:

1. Markiere für jeden Eintrag, ob er **Hassrede** oder **nicht Hassrede** ist.
2. Markiere außerdem Einträge, die Du für **unrealistisch** hältst.

Was meinen wir mit “Hassrede”?

Wir definieren **Hassrede** als *Beleidigungen und Beschimpfungen einer geschützten Gruppe oder Mitgliedern einer geschützten Gruppe auf Grund ihrer Mitgliedschaft in dieser*. Geschützte Gruppen basieren auf Alter, Behindertenstatus, Rasse (Farbe, Nationalität, ethnische oder nationale Herkunft), Religion oder Glaube, Geschlecht, Geschlechtsidentität, oder sexueller Orientierung.

Was meinen wir mit “unrealistisch”?

Wir überlassen das absichtlich Dir. Wir haben keine Einträge mit Absicht unrealistisch gemacht. Diese optionale Markierung soll nur die Datenqualität gewährleisten.

Hilfreiche Tipps:

- **Wörter, die in manchen Kontexten Hass sind, können auch in nicht-Hass gebraucht werden** (z.B. Gegenrede, die hassvolle Schimpfwörter benutzt: *“Wir sollten Leute nicht Neger nennen”*)
- **Das Ziel von Beleidigungen und Beschimpfungen ist entscheidend dafür, ob es sich um Hassrede handelt.** Beschimpfungen von unbestimmten Personen oder nicht-geschützten Gruppen wie Berufsgruppen sind kein Hass (z.B. “Ich hasse Juden” ist Hassrede, aber “Ich hasse dich”, “Ich hasse Pizza” und “Ich hasse Ärzte” ist es nicht).

Andere wichtige Hinweise:

1. Markiere Einträge *nicht* als unrealistisch nur weil sie ungewöhnlich sind. Markiere nur grammatisch inkorrekte Einträge und Nonsens also unrealistisch.
2. Bitte arbeite unabhängig und diskutiere deine Entscheidungen nicht mit Anderen.
3. Die meisten Einträge sind sehr kurz, also zerbrich dir nicht den Kopf an ihnen.

Vielen Dank für Deine Arbeit!