

Resiliente Kommunikation durch Social Media Monitoring, Analyse und Validierung

Abschluss Konsortialtreffen 22. Mai 2024

Rolf Fricke, Jan Thomsen, Ali Sarioglu, Marie Sophie Funk (Condat AG)

Social Media Monitoring zur Gewinnung zusätzlicher, validierter Informationen in Krisensituationen

- Erkennung von Gefahren
- Betroffene Gebiete
- Kommentare, Bilder und Videos von Augenzeugen
- Stimmungslage und Emotionen

Zielgruppen

- Feuerwehr
- Radio/TV-Sender
- Polizei
- THW
- weitere BOS ...

Polizei Berlin @polizeiberlin · Aug 4

Um der @Berliner_Fw die Brandbekämpfung beim Großbrand rund um unseren #Sprengplatz in #Grunewald zu ermöglichen, unterstützen unsere Kolleg. weiterhin.
U.a. sind unsere Wasserwerfer im Einsatz.
Weitreichende Absperungen wurden eingerichtet.
*tsm

21 28 343

Show this thread

Polizei Berlin @polizeiberlin · Aug 4

Das Gelände ist mit Brandmeldeanlagen ausgestattet, verfügt über eine mehrere Meter breite Brandschutzschneise und sieht eine Dauerberegnung der gelagerten #Kampfmittel vor.

35 20 217

Polizei Berlin @polizeiberlin · Aug 5

Wir sind heute mit rund 200 Einsatzkräften vor Ort und danken @PolizeiMV und @PolizeiBB für ihre zur Unterstützung gesandten Wasserwerfer. Gemeinsam konnten wir bisher rund 1,8 Millionen Liter Wasser zur Befeuchtung und Kühlung des Geländes versprühen.

5 4 99

Bundesanstalt THW @THWLeitung · Aug 5

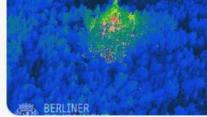
Deutschlandweite Brandeinsätze des #THW: Beim Brand auf dem Munitionsgelände im Berliner #Grunewald unterstützen THW-Einsatzkräfte mit großen Wasser-Pufferbecken, Kraftstofflogistik und Beleuchtung. @BMI_Bund @RegBerlin @Berliner_Fw ...



Berliner Feuerwehr @Berliner_Fw · Aug 4

Update: 15.000 qm Waldfläche brennen. Auch der Sprengplatz Grunewald ist betroffen. Es kommt zu Explosionen. Ein Sperrradius von 1000m wurde festgelegt, die BAB_A_115 ist komplett gesperrt, der angrenzende Bahmverkehr ist komplett eingestellt.
⚠ betreten Sie nicht das Waldgebiet






Tagesspiegel @Tagesspiegel

+EIL+ Großbrand im Berliner #Grunewald: Munitionslagerstätte der Polizei betroffen, Autobahn Avus gesperrt.

Translate Tweet



tagesspiegel.de
Großbrand im Berliner Grunewald: Löscharbeiten auf Sprengplatz beendet – Ur...
Avus wieder offen – Waldgebiet noch gesperrt + Feuerwehr warnt vor Waldbrandgefahren in Berlin + Der Newsblog zum Brand.

Quellen

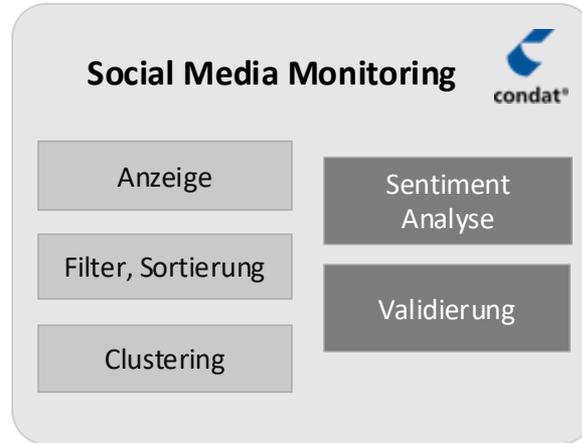


...

kritisches
Ereignis



Suche



weitere
Distributions-
Kanäle
(Broadcast,
Newsroom,
App ...)

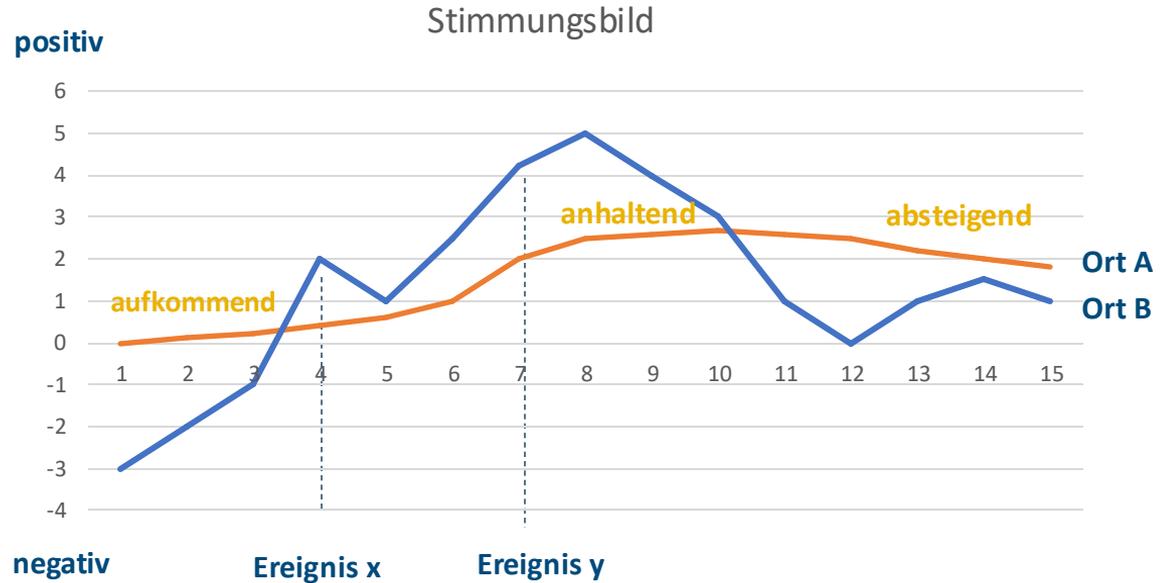


- Leitstelle
- Redaktion



Endnutzer

- Korrelation zu Ereignissen und Orten
- Veränderung der Stimmungslage
- Shitstorm-Ermittlung



Sentiment Analyse

Wörterbuchbasiert

„Feierabend.
Ich **will** nach Hause, Zug, **fällt** aus, das Stellwerk ist kaputt.
Bus auch **schlecht**, da schwerer **Unfall** und alles gesperrt ist.
Wünscht mir **Glück!**“

Gesamtsentiment:

$$\frac{0,2549 + (-0,3529) + (-0,7667) + (-0,7833) + 0,8725}{5} = -0,1551$$

Wort	Sentiment
abbauen	-0,2667
Abendrot	0,8
Fähigkeit	0,5392
fallen	-0,3529
Glück	0,8725
Tisch	0,8725
schlecht	-0,7667
Unfall	-0,7833
wollen	0,2549

Sentiment Analyse - Emotionen

Wörterbuchbasiert

„Feierabend.
Ich will nach Hause, Zug, fällt aus,
das Stellwerk ist kaputt.
Bus auch schlecht, da schwerer
Unfall und alles gesperrt ist.
Wünscht mir Glück!“

Emotionen:

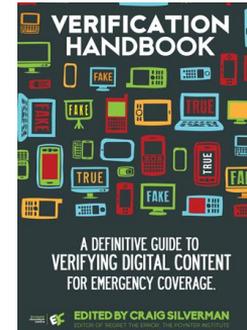
Freude, Furcht, Trauer,
Überraschung, Verachtung

Ekel	Freude	Furcht	Trauer	Überraschung	Verachtung	Wut
Bazille	Applaus	bangen	Abschied	aufschrecken	Abklatsch	Angriff
brechen	betören	grauen	Ende	Entsetzen	charakterlos	Gebrüll
Fäule	Glück	höllisch	freudlos	Knall	Dulli	hassen
igitt	Jubel	Krebs	mies	mitreißen	kaputt	pissig
speien	prima	scheu	Schmerz	radikal	schlecht	sauer
würgen	zujubeln	Unfall	Unfall	Unfall	Trottel	Vorwurf

durch natürlichsprachliche Formulierung der
Anfrage an das OpenAI API:

Antworte mit folgendem JSON Objekt: { "totalSentiment": 0, "emoji": Emoji, "explanationSentiment": [string] } Analysiere die Stimmung im Text mit totalSentiment: Wert für Stimmung von -1 bis 1, emoji: Emoji für Stimmung, explanation: Maximal 3 Stichwörter(Adjektive) auf deutsch für Stimmung. Alle Werte Null wenn kein Text.

- Training für Validierungsprozess durch Best Practice von Journalism Center (Verification Handbook), First Draft News, EU Disinfo Lab ...
- Nutzung von Standardtools für Suche, Mapping Services, Reverse Image Search, ...
- immer Nutzer in der Schleife, um die Ergebnisse der Tools zu bewerten
- Analyse von Quelle, Inhalt, Zeit, Ort, Form, Kontext



von Posts durch ca 130 W3C Credibility Signals:

- Hasssprache, Beleidigungen, Schimpfwörter
- Grammatik- / Rechtschreibfehler
- Emotionalität, Subjektivität
- Ausrufe, Großschreibung
- Vokabularbreite
- Lesbarkeits-Indikator (z.B. Flesch-Kincaid Level)
- Konsistenz von Titel, Zusammenfassung und Inhalt

....



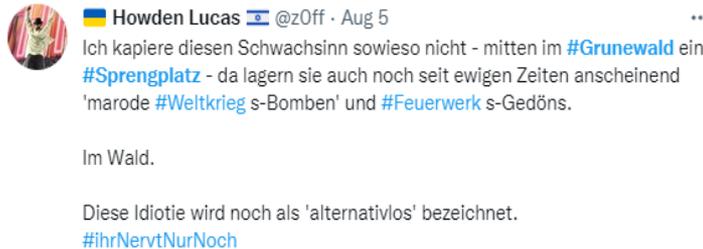
Ziel ist Erhöhung der Transparenz: die KI entscheidet nicht selbst, sondern die Redaktion in der Schleife

EU Regulierungs Maßnahmen

- EU Digital Services Act (DSA) - ratifiziert Feb 2024 für große Plattformen
- EU AI Act – verabschiedet Dez 2023
 - KI-Algorithmen zur Klassifikation, Bewertung (Bias / Diskriminierung), Social Scoring unter Nutzung personenbezogener Daten
- EU Data Act – verabschiedet Nov 2023
 - Offenlegung von Daten gegenüber Behörden im Krisenfall
- EU Market Act – ratifiziert seit Oktober 2022
 - Regelungen zu Gatekeeper Problematik

....

- Entfernung von Hass-Kampagnen, Diskriminierung und illegale Inhalte durch Algorithmen oder Operator
- Wie kann man zwischen kritischer Meinungsäußerung und illegaler Hassrede unterscheiden ? →
- Viele Posts sind nicht wahr oder falsch, legal oder illegal sondern kritische Äusserungen, wie z.B. zu Corona- oder Klimaschutz-Maßnahmen



- Social Networks wie Facebook und Instagram wollen politische Informationen ausblenden
- Twitter Files von Elon Musk: X/Twitter folgt Vorgaben von Regierungen (Türkei Erdogan Wahl, BRD /Corona Politik, Ukraine ...)
- Beiträge mit Wörtern wie z.B. "Hass", "Sklave", "Hetze" werden gelöscht
- Algorithmen und Operator löschen im Zweifelsfall Beiträge um Strafen zu vermeiden
- positiv, da weniger böartige Beleidigungen
- Eliminierung der Kritik von Randgruppen und kritischen Nachfragen
- ➔ Einschränkung der Meinungsvielfalt und Ende des freien Internets ?

- im Krisenfall ist nach DSA ausschließlich EU-Kommission zuständig
- z.B. bei terroristischen Handlungen, Naturkatastrophen wie Erdbeben, Wirbelstürme, Pandemien
- Vorgaben der EU für Anbieter sehr großer Online-Plattformen:
 - Verfahren zur Moderation von Inhalten
 - Vorgaben für Meldungen
 - Anpassung algorithmischer Systeme
- kann sinnvoll sein, aber auch leicht mißbraucht werden

- unstetes Nutzerverhalten bei X/Twitter, TikTok, Facebook ...
- Auf vielen Plattformen nicht alles sichtbar, Untergruppen z.B. nur über bestimmte Accounts (Facebook, Telegram, vk)
- Häufig wechselnde Bedingungen für Kosten und Zugriffsrechte, z.B. bei X/Twitter API
- Zugriff bei X/Twitter nur für verified / bezahlte Accounts und Anzahl der Aufrufe begrenzt
- kontinuierliche Anpassung des Social Media Monitoring Systems



produktreife Version des Social Media Monitorings bis Ende des Jahres

**Vielen Dank für die
Aufmerksamkeit !**

Fragen ?