

## 4 Methoden zur Erfassung, Verwaltung und Auswertung von Tweets

Die in Kapitel 2 vorgestellten Studien verwendeten unterschiedliche Herangehensweisen zur Erhebung und Analyse von Daten auf Twitter: Je nach zeitlicher Fokussierung wurde die Streaming API oder die REST API, in wenigen Fällen auch ein kostenpflichtiger Datensatz eines Datenhändlers verwendet. Zudem wurden die erhobenen Daten in Datenbanken oder einfachen Textdateien gespeichert, die quantitativen oder qualitativen Analysen erfolgten manuell oder computergestützt.

Welche Möglichkeiten der Datenerhebung auf Twitter bestehen und wie diese eingesetzt werden, soll im nächsten Kapitel genauer analysiert werden. In Anlehnung an den üblichen Prozess der Analyse von Online-Daten (Datenerhebung – Datenspeicherung – Datenauswertung) gliedert sich dieser Teil der Arbeit in drei Kapitel. Kapitel 4.1 befasst sich zunächst mit den technischen Aspekten der Datenerhebung und stellt hierfür die beiden kostenlosen Twitter-Schnittstellen gegenüber. Zudem wird auch kurz auf alternative, kostenpflichtige Methoden eingegangen sowie die Vorteile und Grenzen der einzelnen Ansätze anhand praktischer Beispiele dargestellt. Kapitel 4.2 geht schließlich auf den Prozess der Datenspeicherung und -verwaltung ein. Die Wahl einer geeigneten Speichermethode ist entscheidend für die weitere Datenverwaltung, Datenpflege und letztlich auch die Analyse. Neben der einfachen Speicherung in einzelnen Dateien werden unterschiedliche Datenbank-Konzepte angesprochen. Das abschließende Kapitel 4.3 befasst sich mit Analyseverfahren. Für die Auswertung umfangreicher und detaillierter Daten sozialer Netzwerke sind manuelle Verfahren, wie das Kodieren einzelner Tweets, nicht praktikabel, beziehungsweise nur unter großem Aufwand realisierbar. Deswegen steht eine Vielzahl computergestützter Verfahren zur Verfügung, die den Prozess vereinfachen: Von der automatisierten Inhaltsanalyse mit *Bag of Words* Repräsentation bis hin zu vollständigen semantischen Analysen.

## 4.1 Möglichkeiten der Datensammlung

Twitter-Daten können mit unterschiedlichen Methoden erhoben werden. In der Praxis haben sich, je nach Ausgangslage, mehrere Verfahren etabliert. Twitter ermöglicht seit 2014 sechs von 1.300 beworbenen Forschungsprojekten umfassenden Datenzugriff (Krikorian, 2014). Neben diesen *Data Grants* besitzen noch mehrere Institutionen (wie das *MIT* oder die *Library of Congress* in den USA) und Geschäftspartner (z.B. *IBM*, *Brandwatch*) privilegierte Rechte. Die meisten Interessenten für Twitter-Daten müssen jedoch andere, hinsichtlich Umfang und Informationsgehalt deutlich beschränkte, Datenzugänge wählen.

Für „normale“ Twitter-Nutzer besteht die Möglichkeit, eigene Tweets als lokale Kopie zu archivieren, wobei sich der Datenexport aufgrund des Formats und der Einschränkung auf eigene Tweets nicht für Forschungszwecke eignet. Für eine detaillierte Datenabfrage stellt Twitter zwei Typen von Schnittstellen zur Verfügung: die Streaming API und die REST APIs. Diese unterscheiden sich hinsichtlich Datenumfang, Zeitraum und Funktionalität. Eine Nutzung dieser APIs setzt, im Vergleich zu simplen Datenexporten aus Twitter, technisches Know-How und ein entsprechend einsetzbares IT-System voraus. Die in folgenden Kapitel vorgestellten Twitter APIs unterstützen eine Vielzahl gängiger Programmiersprachen zur Steuerung der Datenabfragen und -verarbeitung, wie *Java*, *ASP*, *Ruby* und *Python*.

Alle Schnittstellen benötigen für den Daten-Zugriff einen Twitter-Account und eine darin registrierte Anwendung, die die Datenabfragen verwaltet. Diese App autorisiert schließlich die Endbenutzer für die Verwendung der Twitter APIs, wobei der Zugriff entweder von einem Programm allein (*App Auth*) oder von mehreren Nutzern parallel (*User Auth*) über einen Dienst erfolgen kann. Die Autorisierung muss vor Beginn der Abfrage einmalig durchgeführt werden. Listing 2 zeigt einen typischen Vorgang zur Autorisierung eines Programms bei den Twitter APIs. Hierfür verlangt die API den *Consumer Key* und das *Consumer Secret* aus der App-Verwaltung<sup>8</sup>, die in den Programmcode eingetragen werden. Wenn die Schnittstelle parallel von mehreren Programmen über dieselbe App abgerufen werden soll (also per *User Auth*), benötigt der Anmelde-Prozess zusätzlich *Access Token* und *Access Token Secret*.

---

<sup>8</sup> Der Aufruf der App-Verwaltung erfolgt über <https://apps.twitter.com/app/>

*Listing 2:* OAuth-Autorisierung bei der Twitter API.

```
import tweepy

# Authentifizierung mittels Zugangsdaten, die unter
# https://apps.twitter.com/app abgerufen werden

consumer_key="<Consumer/API Key>"
consumer_secret="<Consumer/API Secret>"
access_key="<Access Token>"
access_secret="<Access Token Secret>"

auth = tweepy.OAuthHandler(consumer_key, consumer_secret)
# Wenn UserAuth-Methode:
auth.set_access_token(access_key, access_secret)

api = tweepy.API(auth)
```

Neben der Möglichkeit, Twitter-Daten mithilfe (selbst) geschriebener Programm-routinen über die Schnittstellen abzurufen, können Daten zusätzlich von Drittanbietern erworben werden. Hierfür werden weder einer Registrierung bei Twitter, noch Programmierkenntnisse vorausgesetzt. Die folgenden Unterkapitel stellen diese drei Varianten kurz vor und vergleichen sie abschließend.

#### 4.1.1 Streaming API

Die Streaming API ist die wohl meistgenutzte Datenquelle, da sie aufgrund ihres Datenumfangs und ihrer Handhabung ideal für großvolumige, langfristige quantitative Analysen ist (Parmelee & Bichard, 2012, S. 56). Twitter-Daten werden hier als konstanter Datenstrom nahezu in Echtzeit gesendet. Man spricht diesbezüglich von einer *push*-basierten API. Es werden also nach einmaliger Anfrage kontinuierlich Daten an den Empfänger übermittelt, wogegen beim *pull*-Prinzip spezifische Daten immer erst nach Aufforderung geliefert werden würden (Kumar et al., 2014, S. 5).

Die Streaming API gliedert sich in drei *Streams*, die einen jeweils unterschiedlichen Endpunkt zur Datennutzung haben. *User Streams* übermitteln alle Daten, die zu einem spezifizierten Nutzer gehören (Tweets der gefolgt Accounts,

Replies auf die eigenen Tweets, und – wenn autorisiert<sup>9</sup> – Direct Messages). Die *Site Streams* bilden die Mehrnutzer-Variante der *User Streams* ab und sind vor allem für öffentliche Webanwendungen gedacht. Während *User* und *Site Streams* nur die Daten einzelner Nutzer enthalten, liefern die *Public Streams* alle öffentlich erhältlichen Daten auf Twitter und ermöglichen, spezifische Nutzer oder Themen zu verfolgen. Sie sind somit am besten zum Data Mining geeignet.

Alle Streams stehen in drei Bandbreiten zur Verfügung: *Spritzer*, *Gardenhose* und *Firehose*, die entsprechend maximal 1, 10 und 100 Prozent aller veröffentlichten Tweets zu einem Suchterm pro Sekunde übermitteln (Gaffney & Puschmann, 2014, S. 57). Die Funktionsweise dieser Limits soll das folgende Beispiel anhand der *Spritzer* genauer erklären: Die APIs benötigen in den meisten Fällen, beziehungsweise für die meisten Abfragemethoden<sup>10</sup> einen Suchterm. Dieser könnte ein Begriff („#merkel“), eine spezifische Nutzer-ID oder ein Koordinaten-Bereich sein. Übersteigt das Gesamtergebnis des Suchterms ein Prozent des momentanen Twitter-Volumens – macht also das potentielle Suchergebnis in einer Sekunde mehr als ein Prozent des gesamten Tweet-Volumens weltweit aus – werden nur ein Prozent der Ergebnisse ausgegeben. Mit welchen Kriterien (z.B. Zeitstempel, Relevanz, Popularität des Nutzers) diese von der API ausgegebenen Tweets gefiltert werden, ist nicht bekannt.

Während *Spritzer* für alle Nutzer kostenlos ist, wird der Zugriff zur *Gardenhose* nur nach Anfrage und in begründeten Fällen (Forschungszwecke, Online-Dienste) erteilt. Der Zugang zur *Firehose* ist stark reglementiert und wird meist nur im Zusammenhang mit wirtschaftlichen, kostenpflichtigen Kooperationen freigegeben: Aufgrund der umfassenden und detaillierten Datenmenge besitzen nur wenige Unternehmen, wie die Datenhändler *Gnip* und *DataSift*, diese Möglichkeit.

Der Datenstrom von *Spritzer* und *Gardenhose* steht außerdem in zwei unterschiedlichen Methoden zur Verfügung: *Sample* und *Filter*. Erstere Variante erzeugt ein Datensample von einem Prozent aller Tweets, wogegen letztere gefilterte übermittelt: Mittels *Track* Befehl können mehrere kommasetrennte Werte abgefragt werden. Die Schnittstelle gibt dann alle Tweets aus, die diese Begriffe beinhalten, sofern das Volumen nicht ein Prozent des momentanen Twitter-Volumens übersteigt. Analog wird *Follow* für Benutzer-IDs und *Locations* für Koordina-

---

<sup>9</sup> Die Autorisierung erfolgt immer auf Nutzerebene. Jeder Anwender kann über die API nur die Direct Messages verwalten, die mit dem jeweilig autorisierten Account verknüpft sind. Forschende haben also keine Möglichkeit, private Nachrichten anderer Nutzer zu lesen.

<sup>10</sup> Ausgenommen der *Sample*-Methoden.

ten verwendet. Ein Sprachfilter steht momentan nicht zur Verfügung. Für alle Methoden gibt es außerdem Einschränkungen bei der Anzahl an Filter-Parameter: So können zeitgleich höchstens 400 Wörter (*Keywords*), 5.000 User-IDs und 25 Orte abgefragt werden (Twitter, Inc., 2015e).

Morstatter, Pfeffer, Liu und Carley (2013) verglichen den Standard-Datenoutput der Streaming API (Spritzer) mit den vollständigen Daten der Firehose hinsichtlich Daten-Abdeckung und Stichprobenqualität. Hierfür wurden unter anderem Tweets nach Hashtag und Ort gefiltert und gegenübergestellt sowie Trend-Themen ermittelt. Die Untersuchung ergab, dass es sich beim 1-prozentigen Sample der *Spritzer* um keine vollständig verlässliche Stichprobe handelt, sondern die Qualität des Samples von den analysierten Begriffen/Themen/Nutzern abhängt. Bei geringen Fallzahlen von Hashtags oder Themen traten kleinere Abweichungen auf. Bei zu großen Fallzahlen, also beispielsweise zu vielen Tweets, die ein bestimmtes Hashtag enthalten, sank die Genauigkeit des Samples. Dies lässt sich damit begründen, dass nach Überschreiten des 1%-Limits der Spritzer nicht mehr alle betreffenden Tweets übermittelt werden. Werden also zu einem Zeitpunkt mehr als ein Prozent aller Nachrichten auf Twitter zu einem gefilterten Begriff /Thema/Nutzer verfasst, wird der Datenstrom auf diesen Anteil (respektive 10% bei der Gardenhose) gedeckelt und die Menge an „verloren gegangenen“ Tweets als Zahl ausgegeben.

Ein etwaiges Erreichen des Limits kann jedoch durch eine vorherige Eingrenzung der Daten mittels mehrerer Filter vermieden werden, sodass letztlich in diesen Fällen alle relevanten Daten vollständig erhoben werden können. Dennoch kann es in einzelnen Fällen, bei sehr großem Tweet-Volumen zu einem einzelnen Thema (beispielsweise bei medialen Großereignissen wie Katastrophen), nicht möglich sein, durch Filter die Tweet-Anzahl so einzugrenzen, um die Bandbreitenbeschränkung nicht zu überschreiten. Auch ist ein sehr spezifisches Filtern der Daten nicht immer erwünscht und sinnvoll (bei globalen oder sehr unspezifischen Themen). In diesen Fällen müssen die Tweets, die durch die Deckelung des Datenstroms nicht übermittelt wurden, nachträglich und manuell mit Hilfe der REST API eingespeist werden (siehe Anwendungsbeispiel in Kapitel 4.1.2). Möglich wären auch aufgeteilte Abfragen, die synchron von unterschiedlichen Endpunkten gestartet werden: Je Hashtag wird eine eigene Abfrage über einen separaten Account mit App gestartet – die Zusammenführung der Daten erfolgt dann während des Sammelns (Schreiben in die gleiche Datenbank) oder nachträglich.

#### 4.1.1.1 Anwendungsbeispiel: Sammeln von Echtzeitdaten auf Twitter

Das folgende Beispiel zeigt einen typischen Prozess zum Sammeln von Tweets: Das Programm in Listing 3 erfasst in Echtzeit Tweets mit Hashtag #Obama und gibt diese direkt im Programmfenster (*Shell*) aus. Zur Vereinfachung des Skriptes wird hier also auf ein Speichern und weiteres Verarbeiten der Tweets verzichtet – dieses Vorgehen wird in Kapitel 4.2 besprochen.

Nach der Autorisierung des Programms bei der Streaming API über die *OAuth*-Methode wird zuerst die Klasse *tweetpylistener* erstellt, die das Python-Paket *Tweepy* nutzt und das Vorgehen bei Eintreffen eines neuen Tweets definiert. Hier werden zur Veranschaulichung des Daten-Outputs alle eingehenden Daten in das Programmfenster geschrieben. Möglich wäre aber auch ein Abspeichern in eine Datei oder Datenbank. Etwaige Fehler erscheinen ebenfalls im Shell-Fenster. Diese Klasse kann jederzeit und an jeder Stelle des Programms aufgerufen werden. Schließlich werden die wesentlichen Parameter der API-Abfrage definiert: Der Suchterm entspricht einer vorher definierten Liste von Begriffen (hier: *obama*). Das Vorgehen bei Tweets, die der Abfrage beziehungsweise dem Suchterm entsprechen, definiert bereits die Klasse *tweetpylistener*, die hier schließlich abgerufen wird. Des Weiteren erfolgt eine explizite Einbindung der Streaming API und die Definition dafür notwendiger Parameter, wie Autorisierungs-Daten und eine Time-Out-Zeit für die Anfragen. Sollte eine Suche nach 600 Sekunden kein Ergebnis liefern, endet der Prozess automatisch. Zuletzt werden noch die Suchfilter über die Variable *list\_terms* integriert.

*Listing 3:* Simpler Vorgang zum Sammeln von Tweets mit dem Begriff „obama“ und direkter Ausgabe im Programmfenster

```
import tweepy

# Authentifizierung
consumer_key="<Consumer/API Key>"
consumer_secret="<Consumer/API Secret>"
access_key="<Access Token>"
access_secret="<Access Token Secret>"
auth = tweepy.OAuthHandler(consumer_key, consumer_secret)
auth.set_access_token(access_key, access_secret)
api = tweepy.API(auth)

# Definiere Vorgehen bei neuen Tweets und Fehlern
class tweetpylistener(tweepy.StreamListener):
```

```

def on_data(self, data):
    print(data)
    return True
def on_error(self, status):
    print(status)

# Definiere Parameter des Streams
if __name__ == '__main__':
    list_terms = ["#obama"]
    listen = tweepylistener(api)
    stream = tweepy.Stream(auth, listen, timeout=600.0)
    stream.filter(track=list_terms)

```

Der Suchfilter track weist eine spezielle Systematik auf, die es erlaubt, Suchbegriffe mit den logischen Operatoren AND und OR zu verknüpfen. Diese Systematik verknüpft Wörter, die nur durch ein Leerzeichen getrennt sind, als Konjunktionen und behandelt kommaseparierte Begriffe als Disjunktionen. Dabei berücksichtigt der Algorithmus des Suchfilters auch explizit Punktationen und Sonderzeichen, wobei diese nicht in #hashtags und @mentions erlaubt sind und somit nur die Suche in normalem Fließtext davon betroffen ist. Tabelle 3 veranschaulicht den Such-Mechanismus.

Tabelle 3: Operatoren des Track Filters der Streaming API. In Anlehnung an Twitter, Inc. (2015f).

| SUCHPARA-METER                | BERÜCKSICHTIGTE BEGRIFFE/TWEETS  | NICHT BERÜCKSICHTIGTE BEGRIFFE/TWEETS    |
|-------------------------------|--|--|
| <b>OBAMA</b>                  | Obama / #Obama / OBAMA / @obama / obama. / http://obama.com                      | BarackObama / #presidentobama / obamaUSA |
| <b>OBAMA'S</b>                | Watching Obama's speech.   | Watching @Obama's speech.                |
| <b>OBAMA SPEECH, OBAMA TV</b> | Watching #obama speech on TV<br>#Obama is on TV!<br>Watching obama speech.       | Nice speech on TV!<br>Watching #obama... |
| <b>OBAMA, CASTRO</b>          | #Obama handshake with #castro<br>#obama is now on TV<br>#castro is meeting POTUS | Obamas handshake with #raulcastro        |

Des Weiteren besteht die Möglichkeit, mittels `follow`, `locations` und `language` Filter nur Tweets bestimmter Nutzer, Orte oder Sprachen zu sammeln. Dieser Mechanismus ist nahezu analog zum `Track` Filter und wird durch kommage-trennte Listen definiert. Logische Operatoren stehen hier allerdings nicht zur Verfügung.

Der Suchfilter aus obiger Streaming API Anfrage ergibt einen beispielhaften Output wie in Listing 4. In Abhängigkeit der Häufigkeit der Suchergebnisse würden sequentiell alle dem Suchfilter entsprechenden Tweets angezeigt werden. Bei einem Fehler erschiene der entsprechende Fehlercode (siehe Code 420 am Ende von Listing 4).

*Listing 4:* Shell-Output der Abfrage aus Listing 3 für einen Tweet. Ausgabe-Code durch Autor gekürzt.

```
{
  "created_at": "Sat Apr 11 11:26:02 +00002015",
  "id": 586852702115663872, "id_str": "58685270211566387",
  "text": "ABC News: Watch President Obama Geek Out While Meeting Usain Bolt . More #Obama #news - http://t.co/OkWn16C9Mn", "source": "\u003ca href=\"http://www.1stheadlines.com\" rel=\"nofollow\"ObamaInTheNews \", \"truncated\": false, \"in_reply_to_status_id\": null, \"in_reply_to_status_id_str\": null, \"in_reply_to_user_id\": null, \"in_reply_to_user_id_str\": null, \"in_reply_to_screen_name\": null, \"user\": { \"id\": 45326213, \"id_str\": \"45326213\", \"name\": \"1stHeadlines\", \"screen_name\": \"ObamaInTheNews\", \"location\": \"USA\", \"url\": \"http://www.1stheadlines.com/obama.htm\", \"description\": \"Breaking news stories about Barack Obama from top online news sources.\", \"protected\": false, \"verified\": false, \"followers_count\": 1616, \"friends_count\": 0, \"listed_count\": 68, \"favourites_count\": 0, \"statuses_count\": 113751, \"created_at\": \"Sun Jun 07 11:48:45 +0000 2009\", [...] \"favorite_count\": 0, \"entities\": { \"hashtags\": [ { \"text\": \"Obama\", \"indices\": [74, 80] }, { \"text\": \"news\", \"indices\": [81, 86] } ], \"trends\": [], \"urls\": [ { \"url\": \"http://t.co/OkWn16C9Mn\", \"expanded_url\": \"http://tinyurl.com/ltvgvk\", \"display_url\": \"tinyurl.com/ltvgvk\", \"indices\": [89, 111] } ], \"user_mentions\": [], \"symbols\": [] }, \"favorited\": false, \"retweeted\": false, \"possibly_sensitive\": false, \"filter_level\": \"low\", \"lang\": \"en\", \"timestamp_ms\": \"1428751562037\" }
420
```



Twitter stellt eine Liste mit Fehlercodes und Erläuterungen zur Verfügung<sup>11</sup>. Die Bedeutung reicht von einfachen Autorisierungs-Fehlern bis hin zu komplexen Problemen wie den *Rate Limits*. Letztere treten auf, wenn gleichzeitig zu viele Anfragen über eine autorisierte App innerhalb eines Zeitfensters gestartet werden. Sollte dieser Fehlercode missachtet werden, droht eine temporäre Sperrung (*Blacklisting*) der IP-Adresse<sup>12</sup>. Um das zu vermeiden, ist eine differenzierte Behandlung einzelner Fehler sinnvoll. Dementsprechend wird die in Listing 3 definierte Klasse `tweepylistener` um einige Fälle erweitert.

*Listing 5:* Erweiterte Suchklasse der Streaming API mit zusätzlichen Fällen

```
class tweepylistener(tweepy.StreamListener):
    def __init__(self, api = None):
        self.api = api or API()
        self.deleted= open('geloescht.txt', 'a')

    # Bestimme Vorgehen für unterschiedliche Tweet-Typen
    def on_data(self, data):
        if "in_reply_to_status" in data:
            self.on_status(data)
        elif "delete" in data:
            delete = json.loads(data)["delete"]["status"]
            if self.on_delete(delete["id"], delete["user_id"]) is
                False:
                return False
        elif "limit" in data:
            if self.on_limit(json.loads(data)["limit"]["track"]) is
                False:
                return False
        elif "warning" in data:
            warning = json.loads(data)["warnings"]
            print warning["message"]
            return False

    # Bestimme Vorgehen in unterschiedlichen Situationen
    # Fall 1: neuer Status-Tweet
    def on_status(self, status):
        print(status)
        return True
```

---

<sup>11</sup> Siehe <https://dev.twitter.com/streaming/overview/connecting>

<sup>12</sup> Weitere Informationen unter: <https://dev.twitter.com/rest/public/rate-limiting>

```

# Fall 2: User löscht Tweet nach gewisser Zeit
def on_delete(self, status_id, user_id):
    self.deleted.write( str(status_id) + "\n")
    return

# Fall 3: Streaming API Rate Limit
def on_limit(self, track):
    sys.stderr.write(time.strftime("%Y%m%d-%H%M%S") +
                    ">> Rate Limit: " + str(track))
    return

# Fall 4: Fehlermeldung mit Fehlercode
def on_error(self, status_code):
    sys.stderr.write(time.strftime("%Y%m%d-%H%M%S") +
                    ">> Fehler: " + str(status_code) + "\n")
    time.sleep(60)
    return False

# Fall 5: Verbindungs-Timeout keine Reaktion
def on_timeout(self):
    sys.stderr.write(time.strftime("%Y%m%d-%H%M%S") +
                    ">> Timeout, warte für 120 Sekunden\n")
    time.sleep(120)
    return

def main():
    list_terms = ["#obama"]
    listener = tweepy.listener(api)
    stream = tweepy.Stream(auth, listener, timeout=600.0)

    while True:
        print time.strftime("%Y%m%d-%H%M%S") +
              ">> Streaming gestartet... beobachte
              und sammle"
        print time.strftime("%Y%m%d-%H%M%S") +
              ">> Suche Twitter nach: " +
              str(list_terms)[1:-1]

        try:
            stream.filter(track=list_terms, async=False)
            break
        except Exception, e:
            time.sleep(60)

if __name__ == "__main__":
    main()

```

Die erweiterte Klasse in Listing 5 definiert nun das Vorgehen des Programms für die unterschiedlichen Arten von Tweets und Fehlermeldungen. Zuerst werden die einzelnen Tweets nach bestimmten Signalwörtern durchsucht und in Typen eingeteilt. Enthält ein Datensatz beispielsweise eine `delete`-Anweisung in Verbindung mit einem Account- oder Status-ID-Datenfeld, handelt es sich um eine Lösch-Anweisung. Enthält er ein `limit`, so handelt es sich um eine Mitteilung, dass das Bandbreiten-Limit überschritten wurde.

Das Vorgehen in den einzelnen Situationen wird schließlich genauer vorgegeben: Bei einem normalen Tweet (inklusive Retweet) wird der gesamte Tweet-Datensatz ausgegeben. Wenn ein `delete`-Befehl enthalten ist, schreibt das Skript die entsprechende Tweet-ID in die vorher definierte Textdatei `geloescht.txt`. So erhält man im Nachhinein eine Auflistung aller gelöschten Tweets. Diese können bei Bedarf manuell entfernt werden. Unter Umständen ist es aber sinnvoll, diese Tweets zu behalten und die Löschanweisung als zusätzliche Information zu verarbeiten: Sei es als Anzeichen späterer Selbstselektion oder als Reaktion auf eine kritische Resonanz durch andere Nutzer.

Im Fall einer Limit-Überschreitung erfolgt die Ausgabe der Meldung im Programm-Fenster. Diese Mitteilung enthält auch die Zahl der nicht erfassten, sozusagen verlorenen Tweets. Bei Meldungen mit Fehlercode wird dieser angezeigt und alle weiteren Abfragen über die Streaming API um eine Minute pausiert. Wenn ein Timeout der Anfrage eintritt, also nach einer vorher definierten Zeit kein dem Filter entsprechender Tweet übermittelt wird, pausiert das Programm für zwei Minuten und setzt danach die Aktivität fort.

Diese Routine behandelt einen Großteil möglicher Probleme automatisch. Einschränkungen, die aufgrund technischer Rahmenbedingungen der Streaming API existieren, wie die Bandbreiten-Limitierung, können hiermit allerdings nicht umgangen werden. In diesem Fall ist es bestenfalls möglich, analog zum Verfahren mit Delete-Anweisungen, die Limit-Meldungen mit der Anzahl der nicht erfassten Tweets in einer separaten Log-Datei zu speichern, sodass zumindest Zahlen über die Missings vorliegen. Ein nachträgliches Erfassen der ausgelassenen Tweets wäre nur über die REST API möglich.

#### 4.1.1.2 Bewertung der Streaming API

Zusammenfassend ist die Streaming API ein geeignetes Mittel zur kontinuierlichen Erfassung einer großen Anzahl von Tweets über einen längeren Zeitraum. Mit Hilfe intelligenter ProgrammROUTINEN läuft der Prozess der Datensammlung nahezu vollständig automatisiert. Die Echtzeit-Daten sind kostenlos und in nahezu

unbegrenztem Umfang verfügbar, wodurch nicht nur Ad-hoc-Analysen, sondern auch Langzeitstudien möglich sind. Zwar ist der Zugriff für die meisten Nutzer und somit auch für viele Forschende nur auf ein Prozent des gesamten Tweet-Aufkommens begrenzt. Eine gezielte Wahl geeigneter Filterparameter könnte dieses Problem in vielen Fällen umgehen. Jedoch besteht dann die Gefahr, relevante Tweets durch zu strikte Filter auszusondern. Hinsichtlich der Zielsetzung, möglichst viele Daten zu sammeln, ist die Verwendung der *Public Streams* die einzige sinnvolle Methode. *User Streams* und *Site Streams* bieten nur eine Plattform für Web-Anwendungen, in denen sich Nutzer mit ihrem Benutzerkonto einwählen können und dienen daher nicht zum Sammeln von Daten.

Zudem gibt es Verzerrungen zwischen dem über die Streaming API zur Verfügung gestellten Sample und der Grundgesamtheit, die momentan noch nicht kompensiert werden können (Morstatter et al., 2013, S. 9). Eine Repräsentativität dieses automatischen Samples ist demnach nicht gegeben. Außerdem gibt es seitens Twitter keine Angabe über die Generierung dieser Stichprobe.

Es gilt zu beachten, dass der Abfrageprozess von Tweets aufgrund der *Push*-Architektur während der kompletten Sammlung laufen muss. Die Streaming API übermittelt kontinuierlich in Echtzeit Daten, die durch ein Skript sofort erfasst werden. Ein nachträgliches Sammeln ist demnach nicht möglich. Somit resultiert aus jeder Unterbrechung der Internetverbindung oder des Prozesses ein Datenverlust.

Die Tatsache, dass kein Abrufen historischer Daten möglich ist, birgt weitere Probleme: Trends müssen frühzeitig (also eigentlich ad hoc) erkannt werden, um möglichst viele relevante Tweets erheben zu können. Dies ist allerdings nur bei langfristig terminierten Ereignissen (wie Wahlen oder Sportereignissen) zuverlässig möglich. Spontane Bewegungen, Themen oder Ereignisse wie politische Skandale oder Naturkatastrophen können mit dieser Methode erst post hoc oder zumindest zeitlich verzögert erfasst werden. Letztlich fehlt bei Echtzeitdaten auch die Möglichkeit, die Anzahl von Retweets oder Favorites direkt abzufragen. Da die übermittelten Twitter-Daten immer Momentaufnahmen zur Veröffentlichung eines Tweets sind, ist die Wahrscheinlichkeit hoch, dass jeder neue Tweet keine Favorites und Retweets hat<sup>13</sup>.

---

<sup>13</sup> Da der Tweet sofort bei Veröffentlichung über die Streaming API übermittelt wird und die Wahrscheinlichkeit sehr gering ist, dass der Tweet Millisekunden nach Veröffentlichung bereits Retweets oder Favorites hat, steht er folglich immer im Ursprungszustand (ohne Retweets und Favorites) im Datensatz.

Eine Erfassung (späterer) Retweets und Favorites wäre innerhalb der *Public Streams* nur über die Betrachtung des jeweils letzten/aktuellsten Retweets möglich: Die Meta-Daten eines Retweets beinhalten immer die Anzahl an Retweets und Favorites des Original-Tweets zum Veröffentlichungszeitpunkt des Retweets (siehe Abbildung 7). Wie die schematische Darstellung zeigt, taucht folglich ein Tweet unter Umständen mehrmals in einer Datensammlung auf: Zumindest einmal als Tweet und eventuell mehrmals als eingebetteter Retweet. Diese Methode könnte letztlich auch Tweets erfassen, die außerhalb des Betrachtungszeitraums liegen, aber als aktuelle Retweets weiterhin im Datenset auftauchen.

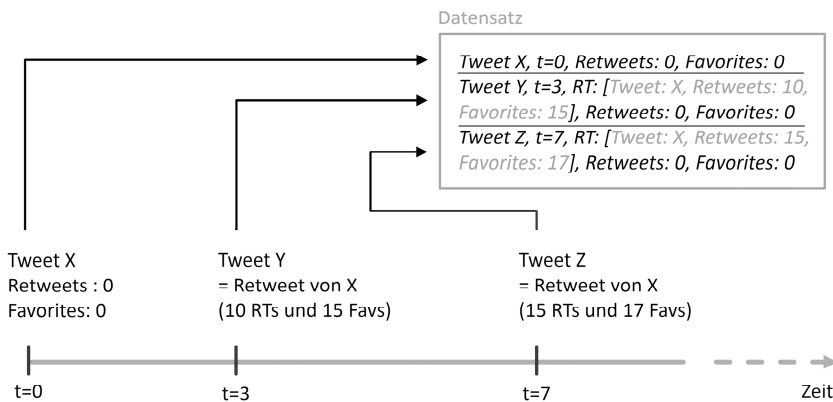


Abbildung 7: Ansatz zum nachträglichen Erfassen von Favorites und Retweets bei der Streaming API. Eigene Darstellung.

Die Datenerfassung könnte also über die Extraktion der eingebetteten Meta-Daten des neuesten Retweets erfolgen. Eine weitere Möglichkeit besteht in der post-hoc Erfassung aller Tweets, wie es die REST APIs erlauben. Da die Daten bei Abruf bereits ein gewisses Alter haben, beinhaltet das nachträglich erhobene Datenset bereits Informationen über Retweets und Favorites bis zum Zeitpunkt der Abfrage. Mit der Funktionsweise sowie den Vor- und Nachteilen der REST APIs befasst sich nun das folgende Kapitel.

### 4.1.2 REST APIs

Während die Streaming API Echtzeitdaten liefert, übermitteln die REST APIs nur vergangenheitsbezogene Daten. *REST* steht für *Representational State Transfer* und ist ein in der digitalen Welt sehr weit verbreitetes Schema der Datenarchitektur und -weitergabe. Es kategorisiert An- und Abfragen in die Operatoren GET, POST, PUT und DELETE. Die REST Schnittstellen bestehen aus einem Bündel an Methoden zur Daten-Interaktion, die sich erheblich von denen der Streaming API unterscheiden. So stehen momentan 35 Methoden für jeweils spezifische Abfragen zur Verfügung: Von den momentan trendigen Themen (*Trending Topics*) bis hin zu den Tweets, Retweets, Blocks, Followern und Favorites eines Nutzers. Über Kombination von Daten mehrerer Abfragen ließen sich zum Beispiel User-Netzwerke oder das Verhalten spezifischer Nutzer/-innen visualisieren. Zudem besteht die Möglichkeit, nicht nur Daten zu lesen (GET-Parameter), sondern auch direkt über die Schnittstelle Tweets zu posten oder Nutzern zu folgen (POST-Parameter) sowie Tweets und Nutzer zu löschen (DELETE).

Für das Sammeln von Tweets eignet sich vor allem die *Search API*, die Teil der REST API ist. Diese funktioniert wie eine Suchmaske und erlaubt Suchanfragen wie: `love OR hate from:mustermann until:2015-01-01` (Suche nach Tweets des Nutzers „mustermann“, die „love“ oder „hate“ beinhalten und vor dem 01.01.2015 verfasst wurden).

Twitter erlaubt in diesem Kontext eine Sortierung der Suchergebnisse: Der Parameter `result_type` ermöglicht eine absteigende Sortierung nach Datum (`recent`), eine Ausgabe absteigend nach Popularität des Tweets (`popular`) oder eine gemischte Ausgabe (`mixed`), welche standardmäßig eingestellt ist (Twitter, Inc., 2015c). Die Sortierung nach `recent` ist für die Datensammlung am praktikabelsten, um eine Vollständigkeit der Daten zu erzielen.

Die REST APIs stellen Daten nicht wie bei der Streaming API über die Push-Methode zur Verfügung, sondern übermitteln diese erst nach einzelnen Pull-Abfragen. Diese Anfragen unterliegen einer starken Reglementierung seitens Twitter. Je nach Abfrage-Methode gibt es unterschiedliche Einschränkungen: Jeder GET oder POST Befehl stellt einen *Request* dar. Sollen beispielsweise Tweets mit einem oder mehreren definierten Keywords abgerufen werden, liefert die API maximal 100 Tweets je Abfrage bei einem Limit von 180 Abfragen je 15 Minuten-Intervall, was ein stündliches Maximum von 72.000 gesammelten Tweets ergibt. Es besteht die Möglichkeit, dieses Limit auf 450 Requests je 15 Minuten zu er-

weitem, wenn die registrierte Twitter-Anwendung nur durch einen Nutzer verwendet wird<sup>14</sup>. Dies ergibt einen Höchstwert von 180.000 Tweets pro Stunde. Ähnlich restriktiv wird auch die Anzahl möglicher Abfrage-Parameter gehandhabt: Die API erlaubt pro Request 20 Keywords (empfohlen wird ein Wert von maximal 10) sowie 100 User-IDs (Twitter, Inc., 2015g). Für die nutzerspezifische Tweet-Suche (Suche über die Nutzer-ID) gilt die Obergrenze von 3200 Tweets. Im Anhang A befindet sich eine Auflistung aller Einschränkungen.

Dieses Beispiel verdeutlicht, dass die Abfrage von Tweets stark eingeschränkt ist. Wächst bei der Streaming API die Obergrenze an erfassbaren Daten je Sekunde noch mit dem gesamten Twitter-Volumen (maximal ein Prozent des Gesamtaufkommens), ist das Limit bei der REST API fixiert. Mehr als 180.000 relevante Tweets können in einer Stunde von einer App nicht gesammelt werden, auch wenn der Anteil am Gesamtvolumen womöglich deutlich geringer als ein Prozent war. Finden sich mehrere Millionen Tweets zu einem Suchterm und sollen diese gesammelt werden, dauert dies unter Umständen mehrere Tage. In diesem Kontext spielt der Verfügbarkeits-Zeitraum von Tweets eine besondere Rolle. Der Zeithorizont für die nachträgliche Suche über die REST APIs liegt momentan bei etwa 6-9 Tagen (Twitter, Inc., 2015g). Meistens besteht keine Möglichkeit, Tweets über diese Periode hinaus abzurufen<sup>15</sup>. Wäre das potentielle Datenset, das durch einen Suchterm angesprochen wird, sehr groß (z.B. bei einem medialen Großereignis wie dem Fußball WM-Finale der Männer 2014 mit über 32 Millionen Tweets), würde die Erfassung über die Search API mehrere Tage dauern. Dies könnte dazu führen, dass das relevante Datenset mit fortlaufender Zeit zum Teil aus dem verfügbaren Zeithorizont rückt und somit nicht mehr erfasst werden kann.

Abbildung 8 veranschaulicht die Problematik in einer schematischen Darstellung. Die Search API sucht in diesem Fall inkrementell nach historischen Tweets zum WM-Finale. Nach jedem Request (zur Erinnerung: ein Request gibt maximal 100 Tweets aus), sucht das Skript nach jeweils vorher veröffentlichten Tweets. Die Suche läuft folglich retrospektiv ab – jeder Request erfasst 100 Tweets, die vor dem vorigen Request datiert sind. Da pro Tag theoretisch maximal 4,32 Millionen Tweets über die Search API gesammelt werden können, würde es über eine Woche dauern, bis alle 32 Millionen relevanten Tweets zum WM-Finale gespei-

---

<sup>14</sup> Da in dieser Arbeit die wissenschaftliche Verwendung der Twitter API betrachtet wird und kein Online-Dienst für mehrere Nutzer, wird davon ausgegangen, dass nur ein Anwender je Twitter-Account/App die Schnittstelle nutzt (App Auth).

<sup>15</sup> Die Begrenzung gilt jedoch nicht für die Tweet-Abfrage auf Nutzer-Ebene über die Timeline (GET statuses/user-timeline). Twitter übermittelt hier maximal die letzten 3200 Tweets eines Nutzers – unabhängig des Veröffentlichungszeitpunktes (Twitter, Inc., 2015d).

chert wären. Da sich der erfassbare Zeithorizont konstant mit der Zeit bewegt (immer die letzten sechs bis neun Tage zum momentanen Abfragezeitpunkt; in der Grafik durch den Zeitraum von  $t=0$  bis  $t=6$  dargestellt), wandern die ersten/frühen Tweets des relevanten Datensets bereits nach einem Tag aus dem erfassbaren Zeitfenster. Einen Tag später (Zeitpunkt  $t=7$  im unteren Bereich der Grafik) könnten zwar wieder bis zu 4,32 Millionen Tweets der letzten sechs bis neun Tage abgefragt werden, jedoch ist der Anfang des relevanten Datensets zu diesem Zeitpunkt bereits außerhalb des verfügbaren Zeitfensters.

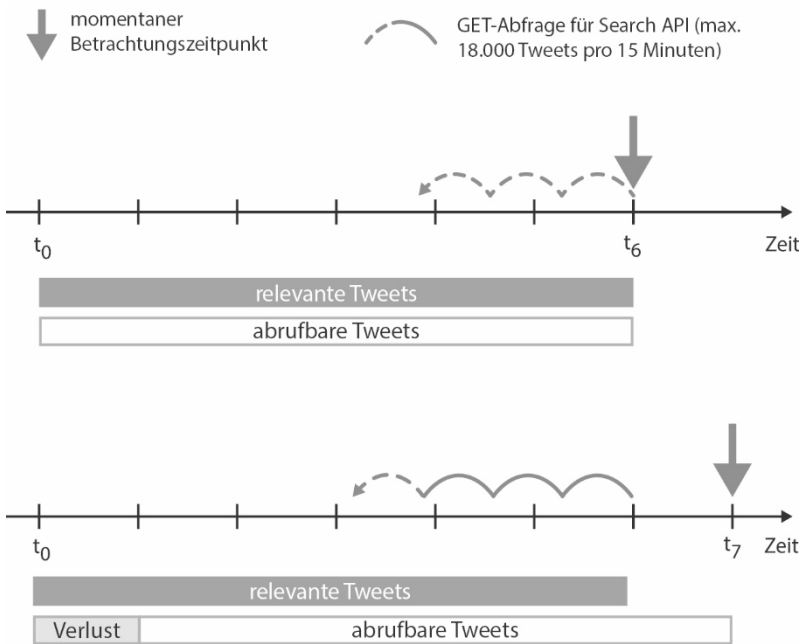


Abbildung 8: Problem der zeitlichen Verfügbarkeit historischer Tweets über die REST APIs. Eigene schematische Darstellung.

Um die negativen Folgen einiger Einschränkungen zu mindern, ermöglicht Twitter bei der REST API eine Eingrenzung der Suchergebnisse über zahlreiche Parameter. Das in Abbildung 8 dargestellte Problem des dynamischen Zeitfensters für



Abfragen ließe sich theoretisch über die zwei Zeitparameter `until` und `since` umgehen: Man startet die Suche am frühestmöglichen Zeitpunkt und sammelt zunächst die Tweets, die als erstes aus dem relevanten und verfügbaren Datenset fallen würden. Die Suche verläuft in diesem Fall entlang des Zeitverlaufs vom ältesten bis zum neuesten relevanten Tweet. Allerdings erlauben diese Parameter nur Datumsangaben und keine Verfeinerung der Zeitgrenzen auf Uhrzeiten. Eine erste Abfrage beliefe sich folglich auf einen ganzen Tag und würde schrittweise alle nachfolgenden Tage erfassen. Dies führt nicht nur zu Problemen mit der strengen Limitierung von Abfragen, sondern ist auch zu ungenau für die Eingrenzung von Suchabfragen.

Relevant für die gängigen Abfrage-Methoden sind deshalb vor allem: `since_id` und `max_id` (Tweet-ID, ab/bis zu der Daten ausgegeben werden sollen). Mithilfe dieser beiden Konstanten sowie `Count` (Höchstzahl an Ergebnissen je Abfrage)<sup>16</sup>, `until` (Stichtag, bis zu diesem Ergebnisse angezeigt werden sollen) und `lang` (Nutzersprache) können beispielsweise Tweets zu einem Hashtag, gestaffelt nach Tweet-ID, abgerufen werden, wodurch sich doppelte Abfrageergebnisse vermeiden lassen und die begrenzten Abfragen ökonomisch sinnvoll genutzt werden können. Jedoch löst auch dieser Ansatz nicht die Problematik des dynamischen Verfügbarkeitshorizonts bei großen Datensets. Sinnvoll wäre also eine Kombination beider Ansätze: Mehrere Skripte mit unterschiedlicher Account- beziehungsweise App-Autorisierung sammeln simultan Tweets. Jeder Account mit eigener App greift dann über die REST APIs Daten eines bestimmten Zeitraums ab (definiert durch `since` und `until`). Durch die gestückelte Datenabfrage reduziert sich dann die Dauer der Datenerhebung.

#### 4.1.2.1 Anwendungsbeispiel: Erheben historischer Tweets

Eine bereits genannte Einschränkung der Streaming API ist der Zeithorizont erfassbarer Tweets: Es können nur Daten ab dem gegenwärtigen Zeitpunkt gesammelt werden. Die Erhebung historischer Tweets kann nur über die REST APIs, beziehungsweise deren integrierte Search API erfolgen. Daneben stehen 35 weitere Methoden zur Verfügung, die jeweils spezifische Daten, wie Followers und Favorites von definierten Nutzern, zur Verfügung stellen. Die hinsichtlich der Datenbreite umfassendste und damit wohl auch meistgenutzte Methode ist jedoch die Nutzung der Search API.

---

<sup>16</sup> Der benutzerdefinierte Wert von `Count` muss unterhalb der von Twitter zugelassenen Zahl an Ergebnissen liegen. Dieses Limit unterscheidet sich nach Abfrage-Methode, liegt jedoch meist bei 100.

*Listing 6:* Einfache Suchabfrage nach Tweets mit „apple“ über die Search API

```
import tweepy

# Authentifizierung
consumer_key="<Consumer/API Key>"
consumer_secret="<Consumer/API Secret>"
access_key="<Access Token>"
access_secret="<Access Token Secret>"
auth = tweepy.OAuthHandler(consumer_key, consumer_secret)
auth.set_access_token(access_key, access_secret)
api = tweepy.API(auth)

#Definiere Suchbegriff
term = "apple"

#Definiere relevante Parameter der Abfrage
for tweet in tweepy.Cursor(api.search,
                            q=term,
                            count=100,
                            lang="de",
                            result_type="recent",
                            include_entities=True).items():
    #Gebe nur Tweet-Zeitpunkt und Text aus
    print tweet.created_at, tweet.text
```

Das Beispiel Listing 6 zeigt eine typische Suche von Tweets, die bis zum gegenwärtigen Zeitpunkt erstellt wurden und ein vorher definiertes Schlagwort enthalten. Die Autorisierung entspricht dem Vorgehen bei der Streaming API. Nach der Definition des Suchterms werden weitere relevante Parameter der Abfrage bestimmt. Zum einen wird die Cursor-Methode verwendet, die die Ergebnisse jedes Requests in ein Cursor-Objekt bündelt. Ein Cursor ist ähnlich wie eine Seite und dient der Aufteilung der gesammelten Daten in Datenblöcke, um diese besser sichten und verarbeiten zu können. Als Ausgabe-Obergrenze für die Suchabfrage wird mit `count` ein Wert von 100 Tweets festgelegt. Dies entspricht dem allgemeinen Limit der REST API für Suchanfragen (siehe Anhang A). Zudem wird das Abfrageergebnis auf deutsche Tweets<sup>17</sup> eingeschränkt und absteigend nach Tweet-Zeitpunkt sortiert. Die Suche integriert auch die Tweet-Entities (also #hashtags, @mentions, URLs etc.), sodass die API auch Tweets ausgibt, die den Begriff nicht

---

<sup>17</sup> Twitter analysiert alle Tweets und erkennt automatisch die geschriebene Sprache. Jedoch ist dieser Spracherkennungs-Mechanismus nicht vollkommen zuverlässig. Siehe hierzu auch Kapitel 4.3.1.

im Tweet-Text, aber in einer URL<sup>18</sup> enthalten. Zur Vereinfachung der Darstellung wird der Output schließlich auf den Tweet-Zeitpunkt und die Nachricht beschränkt.

*Listing 7:* Shell-Output für Programmcode aus Listing 6

```
>>>
2015-04-13 13:12:27 RT @onlinekosten: Apple Watch: Fast eine
Million Vorbestellungen in den USA. Welche Modelle sind
besonders gefragt? http://t.co/Hw5eecA6x5 ...
2015-04-13 13:12:11 New Free iPad app - Apotheke im Kaufland
Lörrach - http://t.co/GHc5bR52uh
2015-04-13 13:12:10 #itunes #iphone5 PhotoStitcher - Maxim
Gapchenko http://t.co/QEjDVJMYu3 #apps #apple
2015-04-13 13:11:31 @timohetzel das Video ist komplett
nervig, ansonsten ist der Stick aber überraschend gut.
Dagegen sieht der Apple TV (bisher) alt aus.
2015-04-13 13:11:23 RT @DJUNDERGROUND: #NowPlaying "Player
Party-IceBerg x @JTMoneyMIATL x @YoungTrizo" on @AL-
LOUTHUSTLE.COM Listen http://t.co/10GmimxmzJ #Clu...
2015-04-13 13:10:51 @CHIP_online #Apple hats vorgemacht
```

Listing 7 stellt beispielhaft das Ergebnis einer solchen Abfrage dar. Dabei zeigt sich, dass der Suchterm nicht zwangsläufig im Tweet-Text vorkommen muss: Der zweite und fünfte Tweet enthalten nur in den Meta-Daten den Begriff „apple“ (hier im Link, der durch Twitter automatisch zu einem *t.co*-Shortlink umgewandelt wird) und werden deshalb von der Search API als relevantes Ergebnis angesehen. Bei der Formulierung des Suchterms ist außerdem zu beachten, dass sich nicht nur der Mechanismus, sondern auch die Systematik der Suchterm-Formulierung von der Streaming API unterscheidet. Eine durch Komma getrennte Liste mehrerer Begriffe ist nicht möglich, der logische Operator OR wird direkt zwischen zwei Suchbegriffe geschrieben. Tabelle 4 stellt die weitere Systematik dar.

---

<sup>18</sup> Twitter kürzt alle URLs und wandelt diese in *t.co*-Links um. So wird aus [www.link.de/ziel](http://www.link.de/ziel) etwas wie [t.co/abcde123456](http://t.co/abcde123456). Da Twitter die ursprüngliche URL in den Meta-Daten einbettet, können diese Teile der URL als Suchergebnis dienen.

Tabelle 4: Such-Operatoren der Search API. In Anlehnung an Twitter, Inc. (2015h).

| SUCHPARAMETER                       | BERÜCKSICHTIGTE BEGRIFFE/TWEETS   |
|-------------------------------------|---|
| <b>apple tv</b>                     | Tweets mit beiden Begriffen (Position egal)                               |
| <b>“apple tv”</b>                   | Tweets mit der genauen Wortfolge  |
| <b>apple or tv</b>                  | Tweets mit einem dieser Begriffe  |
| <b>apple -tv</b>                    | Tweets mit <i>apple</i> , aber ohne <i>tv</i>                             |
| <b>#apple</b>                       | Tweets mit Hashtag <i>apple</i>   |
| <b>from:userx</b>                   | Tweets des Users <i>userx</i>   |
| <b>to:userx</b>                     | Tweets an User <i>userx</i> (Replies oder Mentions)                       |
| <b>@userx</b>                       | Tweets mit Mention von <i>userx</i>                                       |
| <b>apple or tv since:2015-07-01</b> | Alle Tweets mit Begriffen <i>apple</i> oder <i>tv</i> seit dem 01.07.2015 |
| <b>apple until:2015-07-01</b>       | Alle Tweets mit Begriff <i>apple</i> bis zum 01.07.2015                   |
| <b>apple :)</b>                     | Alle Tweets mit Begriff <i>apple</i> und positiver Stimmung               |
| <b>tv?</b>                          | Alle Tweets mit Begriff <i>tv</i> und gestellter Frage                    |
| <b>apple filter:links</b>           | Alle Tweets mit Begriff <i>apple</i> und einer enthaltenen Verlinkung     |

Die gezeigte Abfrage über die Search API ist zum Beispiel als Vorbereitung für eine umfassende und langfristige Erhebung mittels der Streaming API, bei der zunächst der aktuelle Datenbestand gesichtet werden soll, oder für eine rein historische Datensammlung denkbar. Möglich wäre aber auch eine ergänzende Abfrage von Tweets, die aufgrund eines Überschreitens der Bandbreiten-Höchstgrenze nicht mehr erfasst werden konnten. Dies ist wiederum mit Hilfe ergänzender Suchoperatoren möglich. Hierfür stehen die Parameter `since_id` und `max_id` zur Verfügung, die die Such-Ergebnisse mit einer unteren und oberen Tweet-ID-Grenze einschränken können. Da Tweet-IDs fortlaufend nummeriert werden, können – bei Kenntnis einer bestimmten Tweet-ID – alle Tweets bis oder ab diesem Wert angefordert werden.

Es gilt jedoch zu beachten, dass die REST APIs die Anzahl von Abfragen in einem bestimmten Zeitintervall restriktiv behandeln (siehe oben). Diese Problematik sollte die Programmierung automatisierter Abfragen berücksichtigen, indem zum einen die Suchanfragen in Sequenzen unterteilt werden und zum anderen die Zahl der bereits getätigten Requests und respektive die Zahl der noch möglichen Abfragen kontrolliert wird. Tweepy ermöglicht sowohl eine Eingrenzung der Abfragen mithilfe der Parameter `since_id` und `max_id`, als auch eine Überwachung der Abfrage-Limitierung.

In Listing 8 wird der Code um die entsprechenden Funktionen erweitert. Das Programm soll in diesem Fall nachträglich bis zu eine Million Tweets mit dem Begriff „tatort“ erfassen. Zudem speichert die Suchschleife alle ermittelten Tweets in einer Textdatei. Das Programm sucht in sequentiellen Abfragen retrospektiv nach Tweets, die dem Suchterm entsprechen. Dabei findet auch ein Abgleich nach Duplikaten statt. Der Suchbereich wird anhand automatisch gesetzter Werte für `sinceID` und `maxID` dynamisch begrenzt: die Schranken verschieben sich je Anfrage, wobei in diesem Fall die obere Grenze der ID des jeweils zuletzt gespeicherten Tweets entspricht und die untere Grenze nicht definiert ist. Möglich ist aber auch eine manuelle Festlegung auf bestimmte IDs. Pro 15-Minuten-Intervall tätigt das Skript 450 Requests, die jeweils maximal 100 Tweets sammeln. Diese Werte entsprechen den maximal zugelassenen Interaktionen per AppAuth über die REST API, welche während der gesamten Programmaktivität durch das Programm überwacht werden. Dies gewährleistet, dass die Limits nicht überschritten werden, sondern das Programm bis zum nächsten Zeitfenster pausiert. Die Zahl der gesammelten Tweets je Sequenz wird addiert und mit einer zuvor definierten Obergrenze abgeglichen. Bei Erreichen dieses Limits endet der Prozess.

Somit ist diese Programmschleife zum nachträglichen Sammeln von Tweets geeignet: Nach Initiierung erfolgt der Erfassungs- und Speicherprozess selbstständig, in Abhängigkeit der verfügbaren Requests. Es gilt jedoch zu beachten, dass eine Änderung der Parameter nach Start des Prozesses nicht mehr möglich ist. Bei einer Anpassung des Suchterms oder der oberen/unteren Tweet-ID muss der Prozess neu gestartet werden.



```
tweetCount += len(newTweets)
print(time.strftime("%Y%m%d-%H%M%S") + " | {0} Tweets
      gedownloaded".format(tweetCount))
maxID = newTweets[-1].id
except tweepy.TweepError as error:
print("Fehler : " + str(error))
break
print("-----")
print("{0} Tweets in {1} gespeichert".format(tweetCount,
      exportfile))
```

#### 4.1.2.2 Bewertung der REST APIs

Grundsätzlich sind die REST APIs ein umfassendes Toolkit mit einer Vielzahl an Einsatzmöglichkeiten und damit ähnlich geeignet für das Sammeln von Daten, wie die Streaming API. Allerdings hat die Schnittstelle eine grundlegend andere Funktionsweise als die zuvor betrachtete Streaming API. Es werden nur historische Daten übermittelt und diese nur nach spezifischer Anfrage. In ihrer Funktionsweise ist die Schnittstelle sowohl zur Einarbeitung beziehungsweise Sichtung der Datengrundlage zu Beginn eines Forschungsprojektes, als auch zur Erhebung vollständiger Datensätze ideal. Denkbar ist auch eine Verwendung als zusätzliche Datenquelle, falls das Bandbreiten-Limit der Streaming API überschritten wurde und diese nicht erfassten Tweets nachträglich in den Datensatz aufgenommen werden sollen. Des Weiteren stehen umfassende Abfragemöglichkeiten zur Verfügung: Neben trendigen Themen (*Trending Topics*) oder an einem bestimmten Ort veröffentlichte Tweets auch sehr spezifische Informationen, wie die Friends eines Twitter-Users oder dessen Favorites und Retweets. Dies erlaubt nicht nur das Abbilden des historischen Nutzerverhaltens, sondern auch ganzer Nutzer-Netzwerke.

Andererseits bewirkt die Fülle an Abfragen und Kombinationsmöglichkeiten auch eine höhere Komplexität der Datenverarbeitung. Es stehen insgesamt 36 unterschiedliche Methoden zur Verfügung, die sich sowohl hinsichtlich ihrer Funktion, als auch ihrer Parameter unterscheiden. Um verwertbare Daten zu erhalten, sind oftmals mehrere Abfragen mit unterschiedlichen Methoden notwendig, deren Ergebnisse dann verknüpft werden müssen. Um beispielsweise das Netzwerk eines Nutzers aufzulisten, sind zwei separate Requests notwendig: GET friends/ids und GET followers/ids, die dann wiederum mit GET statuses/users/ids abgeglichen werden müssen. Des Weiteren ist die Anzahl an Requests streng reglementiert. Je nach Methode dürfen innerhalb eines 15-Minuten-Intervalls 15 bis 450 Abfragen gestartet werden, die wiederum einen limitierten

Datenumfang haben. Dadurch wird zwar eine Überlastung der Twitter-Server vermieden, allerdings erhöht die unübersichtliche Struktur an Einschränkungen ebenfalls die Komplexität von Abfragen. Nutzer der REST API müssen immer darauf achten, die jeweiligen Höchstgrenzen einzuhalten, um nicht eine Sperrung der IP-Adresse zu riskieren.

Wie bereits erwähnt, stellt die REST API nur historische Daten zur Verfügung. Dies hat den Vorteil, wichtige Daten ex post abrufen zu können. Da sich Trends beziehungsweise populäre Themen auf Twitter meist sehr schnell formieren und ebenso schnell wieder abebben, ist eine nachträgliche Daten-Erhebung für wissenschaftliche Zwecke sehr praktikabel. Somit muss man nicht innerhalb kürzester Zeit auf Trends reagiert oder einen etwaigen Verlust von Tweets vor Betrachtung eines Themas befürchten. Allerdings ist der Zeitraum für die vergangenheitsbezogene Tweet-Erfassung ebenfalls eingeschränkt: Laut Twitter (2015g) stehen über die Search API nur Tweets zur Verfügung, die nicht älter als 6 bis 9 Tage sind. Somit ist beispielsweise eine nachträgliche Betrachtung von älteren Themen/Ereignissen über diese Schnittstelle nicht möglich. Zudem gilt es zu beachten, dass bei sehr bedeutenden Ereignissen mit mehreren Millionen potentiell relevanten Tweets die Gefahr besteht, aufgrund der sehr eingeschränkten Sammel-Kapazität von 180.000 Tweets pro Stunde nicht rechtzeitig alle Tweets zu sammeln, bevor diese aus dem verfügbaren Zeitfenster verschwinden (siehe weiter oben).

Letztlich sollte auch die Qualität des Datensets kritisch betrachtet werden: Twitter (2015g) gibt an, dass die REST API Daten nicht mit dem Anspruch auf Vollständigkeit übermittelt, sondern die Ergebnisse nach Relevanz auswählt. Die Zahl der verfügbaren historischen Tweets hängt von der Relevanz des Themas im gesamten Tweet-Volumen ab. Bei unbedeutenden Themenkomplexen kann der Zeithorizont auch geringer als die mögliche Zeitspanne von 6 bis 9 Tagen ausfallen. Beziehungsweise kann es durchaus sein, dass nicht alle Daten, die den Filterkriterien entsprechen, zur Verfügung stehen.

Das Konzept der Vollständigkeit der Twitter-Daten verfolgt vor allem der kostenpflichtige Anbieter *Gnip*. Die Möglichkeit, Tweets und deren Metadaten über Drittanbieter zu sammeln oder über Datenhändler zu erwerben, soll im Folgenden skizziert werden.



### 4.1.3 *Drittanbieter*

Zusätzlich zu der manuellen Datenerhebung über die Twitter APIs stehen spezialisierte Online-Dienstleister zur Verfügung. Die bekanntesten sind *Gnip*<sup>19</sup> und *DataSift*<sup>20</sup>. Im April 2015 gab jedoch Datasift bekannt, dass das Unternehmen ab August 2015 den Zugriff auf Twitter-Daten verliert (Halstead, 2015). In Zukunft soll Gnip der einzige Anbieter für Twitter-Daten sein, womit sich Twitter die alleinige Kontrolle über die kommerzielle Verwertung seiner Daten schafft (Hofer-Shall, 2015). Auf Basis einer monatlichen Gebühr (mehrere tausend USD) können unter anderem repräsentative Stichproben (z.B. das 10%-Sample *Decahose* von Gnip) oder historische Tweets bis zu einem Stichtag angefordert werden. Das verfügbare Datenset ist somit hinsichtlich der Informationstiefe sehr ausführlich<sup>21</sup>, weist jedoch im Normalfall in der Breite nicht den Datenumfang auf, der theoretisch mit gebündelten Abfragen über beide APIs abgerufen werden könnte. Gnip beschränkt sich hier, aufgrund der großen Datenmenge, auf die gängigsten Metadaten (Gaffney & Puschmann, 2014, S. 58). Statt einem Request über einer der APIs müssen hier spezifische Datenanfragen an den Dienstleister gestellt werden, der dann – je nach Datentyp – die relevanten Daten einmalig oder regelmäßig übermittelt. Die Daten sind bereits strukturiert und können beispielsweise als Datenbank-Import oder mit vorher extrahierten, relevanten Werten versendet werden.

Neben Gnip steht eine größere Zahl kostenloser oder gebührenpflichtiger Programme oder Dienste zur Verfügung, die das Sammeln von Tweets vereinfachen, jedoch teilweise nur eingeschränkte Funktionen beinhalten. Als beispielhafte Auswahl dieser Dienste stellt diese Arbeit die zwei gängigen Online-Dienste *Twitonomy* und *Tweet Archivist* vor.

Twitonomy (Diginomy Pty Ltd., 2015) erlaubt in seiner kostenlosen Version unter anderem das Sammeln von Tweets, ein Aggregieren von Tweets eines Nutzers sowie grundlegende Analysen über dessen Twitter-Verhalten, Interaktion und Vernetzung. Die kostenpflichtige Variante beinhaltet detailliertere Informationen und eine Funktion zum Speichern beziehungsweise Exportieren von Tweets und Analysen als PDF oder Excel-Sheet. Twitonomy ist sehr einstiegfreundlich und gibt mit geringem Aufwand viele aufschlussreiche und übersichtlich dargestellte Informationen. Ebenso ist der Export der Daten intuitiv. Dennoch eignet sich der Dienst nur für die Einarbeitung in die Twitter-Forschung, oder – bei geringem Informationsbedarf – für sehr grundlegende Analysen. Die Datenverfügbarkeit ist

---

<sup>19</sup> *Gnip* wurde 2014 von Twitter übernommen. [www.gnip.com](http://www.gnip.com)

<sup>20</sup> [www.datasift.com](http://www.datasift.com)

<sup>21</sup> Bei *Gnip* sind alle Tweets seit dem Tag der Freischaltung des Online-Dienstes (21.05.2006) gespeichert.

sehr eingeschränkt: Es stehen nur die etwa 3000 neuesten Tweets zu einem Hash-tag oder Nutzer zur Verfügung. Zudem übermittelt Twitonomy nur einen sehr rudimentären Datenumfang: Neben Datum, Nutzernamen und Tweet-Text umfasst ein Datensatz nur den Link zu diesem Tweet, die Client, über den der Tweet geschrieben wurde sowie Angaben über Followers, Retweets und Favorites. Das Datenmaterial ist dabei immer historisch. Twitonomy verwendet also die REST APIs.

Demgegenüber ermöglicht Tweet Archivist sowohl retrospektive Abfragen (bis zu maximal 2000 Tweets) also auch das Einrichten von zukünftigen Abfragen über die Search API (Tweet Archivist, Inc., 2015a). Die Kernfunktion ist das Definieren von Suchtermen, für die der Dienst (stündlich aktualisiert) alle zugehörigen Tweets ausgibt. Die Ausführung der REST API-Abfragen wird folglich auf die Server von Tweet Archivist ausgelagert, welcher die Suchabfragen koordiniert. Zusätzlich besteht auch hier die Möglichkeit, grundlegende Statistiken über Nutzer, Inhalte oder die zeitliche Verteilung abzurufen. Momentan können bei einer monatlichen Pauschale bis zu drei Archive, also drei Suchabfragen, eingerichtet werden.

Dennoch eignet sich auch Tweet Archivist nur für die Einarbeitung in die Twitter-Forschung oder für sehr einfache und limitierte Abfragen. Das Unternehmen gibt eine maximale Archivgröße von 50.000 Tweets an, danach erfolgt eine neue Befüllung mit Daten (Tweet Archivist, Inc., 2015b).

Zusammenfassend sind beide vorgestellten Online-Dienste nur eingeschränkt für die wissenschaftliche Nutzung zu empfehlen. Aufgrund des begrenzten Datenumfangs (hinsichtlich Breite und Tiefe) besteht kaum eine Möglichkeit, ausreichend große Datensets zu erzeugen. Zudem fehlt ein Einblick in die Erzeugung des übermittelten Datensets: Die Daten werden aus einer *Black Box* geliefert, ohne dass Einblicke in die Auswahlmethoden oder die Zuverlässigkeit der Datenerfassung (Missings, Latenzen) gewährt werden (Gaffney & Puschmann, 2014). Gerade im Hinblick auf Reliabilität sind Online-Dienste somit für die Erhebung von Daten für wissenschaftliche Analysen ungeeignet, sondern eher für das Einarbeiten in die Thematik. Demgegenüber ermöglicht Gnip einen umfassenden Zugriff auf theoretisch alle verfügbaren Twitter-Daten. Jedoch sind alle Datenanfragen mit sehr hohen Kosten verbunden.

#### 4.1.4 Vergleich der Möglichkeiten zur Datensammlung

Ein Vergleich der vorgestellten Möglichkeiten zur Datenabfrage ist schwierig, da sich alle Varianten in Zeitraum, Umfang und Methode der Erhebung grundsätzlich unterscheiden (siehe Abbildung 9). Die Streaming API übermittelt nach einmaliger Initiierung des Abfrageprozesses zeitlich unbegrenzt Echtzeitdaten, die anhand von Filtern eingegrenzt werden. Dagegen erlauben die REST APIs nur die Erfassung historischer Daten eines begrenzten Zeitraums, wobei diese explizit (und bei höherem Datenumfang mehrfach und sequentiell) angefragt werden müssen. Der Datenanbieter Gnip bietet wiederum Daten von sowohl vergangener als auch zukünftiger Tweets, jedoch in einem eingeschränkten Informationsgrad und vor allem bei einer hohen monatlichen Gebühr. Gerade hinsichtlich des finanziellen Aspektes ist die Variante, Daten über Drittanbieter zu erhalten, wohl für viele Forschungszwecke unattraktiv.

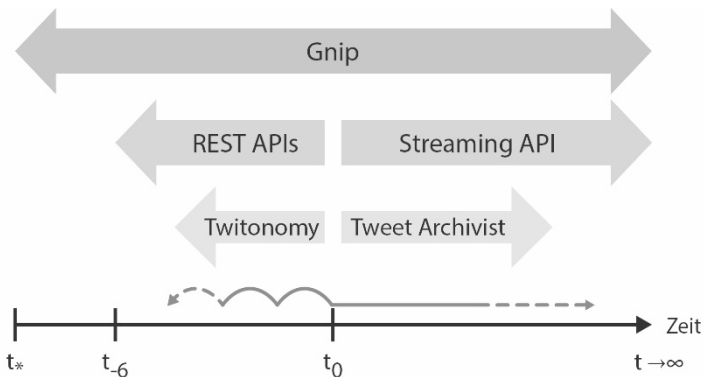


Abbildung 9: Zeithorizont verschiedener Methoden zur Datensammlung. Eigene Darstellung.

Die Streaming API bietet ideale Möglichkeiten zum Sammeln eines möglichst kompletten und zeitlich unbegrenzten Datensatzes relevanter Tweets (oder eines Samples aller Tweets) ab einem definierten Zeitpunkt, wogegen hier die Möglichkeit fehlt, historische Tweets zu erfassen. Sofern die Menge relevanter Tweets nicht ein Prozent des Gesamtvolumens überschreitet, werden alle Daten ohne Einschränkung übermittelt. Sinnvoll gesetzte Abfrage-Filter können das Risiko einer Bandbreitenüberschreitung auf ein akzeptables Niveau senken. Jedoch reichen

auch diese Filter bei sehr populären Themen (wie Katastrophen) unter Umständen nicht aus, um das Suchergebnis ausreichend einzugrenzen. Wie bereits aufgezeigt, besteht aufgrund der Funktionalität der Schnittstelle keine Möglichkeit, Tweets vor dem gegenwärtigen Zeitpunkt zu sammeln – diese Funktion ersetzen die REST APIs.

Die REST APIs beziehungsweise deren Search API eignen sich hervorragend für frühe Phasen eines Forschungsprojektes: Die Abfrage historischer Twitter-Daten ermöglicht einen Überblick der Datenstruktur, Kommunikationsweisen und Häufigkeiten. Die Begrenzung der Datenabfrage hinsichtlich Umfang und Zeitintervall kann mit geringem Programmieraufwand kontrolliert und durch sequentielle Abfragen mit entsprechenden Parametern nahezu umgangen werden (siehe Listing 8). Denkbar und sinnvoll ist auch die nachträgliche Suche von Tweets, die bei Überschreiten des Bandbreiten Limits der Streaming API nicht erfasst wurden. Dies bedarf allerdings einer Kenntnis der Zeitpunkte, zu denen eine Überschreitung stattfand, um über die erste und letzte Tweet-ID dieses Zeitfensters die Suche nach Tweets einzugrenzen. Deshalb empfiehlt sich eine Speicherung der Limits in Logfiles. Ein anschließender Dubletten-Abgleich filtert schließlich bereits vorhandene Tweets. Jedoch ist der Zeitraum für die Datenerfassung über die REST API nur auf etwa eine Woche begrenzt und eine Vollständigkeit des Datensatzes für diese Zeitspanne nicht garantiert. Dementsprechend sollte die Nutzung der REST APIs immer mit Vorsicht erfolgen.

Daher eignen sich beide Twitter-Schnittstellen nicht zur Generierung valider Daten aus der größeren Vergangenheit für Forschungszwecke. Hierfür müssen Informationen vom Drittanbieter Gnip erworben werden. Während bei der Streaming API das Datenvolumen unter Umständen zu groß ist und durch Filter beschränkt werden muss, besteht bei der REST API wiederum die Gefahr, dass nicht (mehr) alle Daten übermittelt werden können. Dieses Problem gibt es bei gekauften Daten nicht: Gnip besitzt einen exklusiven und uneingeschränkten Zugriff auf alle Tweet- und Nutzer-Daten und speichert diese für mehrere Jahre. Dies ermöglicht auch Langzeitstudien oder komparative Analysen der Daten verschiedener Jahre. Die Daten sind jedoch, in Abhängigkeit vom gewünschten Umfang, kostspielig.

Sowohl über die Streaming API als auch über die REST API lassen sich weder valide Vollerhebungen, noch verlässliche Zufallsstichproben generieren. Die Ausgabe der REST APIs basiert, wie bereits angemerkt, auf Relevanz und nicht Vollständigkeit. Selbst wenn über die Streaming API das Limit von einem Prozent nicht überschritten wird, gewährleistet Twitter keine vollständige Übermittlung

aller relevanten Tweets zum Suchterm. Die Betrachtung birgt bei diesen beiden Datenquellen immer ein Risiko.

Bei Daten, die durch die Twitter APIs gesammelt wurden, besteht im Hinblick auf rein quantitative Analysen zudem die Gefahr der Fehlinterpretation von Tweet-Volumina. Aufgrund des unbekanntes Gesamtvolumens auf Twitter könnten eigentlich natürliche Veränderungen im Tweet-Volumen überbewertet werden. Gerade bei geringen Betrachtungszeiträumen sind Daten anfällig für tages-, zeit- und nutzerbedingte Schwankungen, beispielsweise durch Wetter, Tag-Nacht-Wechsel, Zeitzonen oder dem individuellen Nutzerverhalten. Eine plötzliche Spitze in der Hashtag-Häufigkeit bedeutet demnach nicht zwangsläufig eine hohe Popularität, sondern kann auch aus einem allgemeinen Anstieg der momentan aktiven Nutzerzahl (aufgrund der Tageszeit) resultieren. Demnach sollte eine rein quantitative Betrachtung von Tweet-Häufigkeiten immer im Kontext der Gesamtzahl aller Tweets zu diesem Zeitpunkt betrachtet werden. Es besteht allerdings nur die Möglichkeit, diese Werte über die eingeschränkt verfügbare, beziehungsweise kostenpflichtige *Firehose* abzurufen (bei Gnip). Eine Einordnung der erfassten Werte ist somit in den meisten Fällen mit kostenlosen Mitteln kaum oder gar nicht möglich. Tabelle 5 (nächste Seite) fasst die jeweiligen Vor- und Nachteile der vorgestellten Quellen nochmals zusammen.

Sollten Tweets ohne zusätzliche Kosten erhoben werden, ist letztlich die Art der zu erhebenden Daten für die Wahl der Datenquelle entscheidend. Hinsichtlich der Datenstruktur sind beide Schnittstellen nahezu identisch. In manchen Fällen ist eine Kombination beider APIs die praktikabelste Methode zur Datenerfassung auf Twitter, um die jeweiligen Beschränkungen zu mindern, beziehungsweise Vorteile bestmöglich zu nutzen. Die REST APIs dienen dabei zur Sichtung, Einarbeitung sowie zur Erhebung vergangener Tweets oder nutzerspezifischer Informationen. Außerdem werden sie, bei Überschreiten des Bandbreitenlimits der Streaming API, zur nachträglichen Erfassung „verloren gegangener“ Tweets eingesetzt. Die Nutzung der Streaming API erfolgt dagegen ab einem definierten Zeitpunkt zur Sammlung aller zukünftigen Tweets. Somit werden die Stärken beider APIs kombiniert: Die mächtigen Suchwerkzeuge der REST API und deren Fähigkeit, zeitgleich eine große Anzahl an Parametern abzufragen, und die umfangreichen Datenströme der Streaming API, die in Echtzeit eine zeitlich unbegrenzte Zahl von Tweets inklusive aller Meta-Daten übermitteln.

Nachdem nun einzelne, gängige Ansätze zur Datensammlung präsentiert sowie die Vor- und Nachteile für die Forschung aufgezeigt wurden, folgt im nächsten Kapitel eine Betrachtung von Möglichkeiten zur Speicherung und Verwaltung der gesammelten Daten.

Tabelle 5: Vergleich der Quellen für Twitter-Daten

|                      | STREAMING API   | REST APIs  | DRITTANBIETER  |
|----------------------|---|--|--|
| <b>EIGENSCHAFTEN</b> | <ul style="list-style-type: none"> <li>• Echtzeit-Daten</li> <li>• 3 Streams</li> <li>• Daten in JSON-Format</li> </ul>   | <ul style="list-style-type: none"> <li>• Historische Daten</li> <li>• 36 Abfrage-Methoden</li> <li>• Daten in JSON-Format</li> </ul>   | <ul style="list-style-type: none"> <li>• Echtzeit- und historische Daten (unterschiedlich)</li> <li>• Mehrere Datenformate möglich</li> </ul>  |
| <b>VORTEILE</b>      | <ul style="list-style-type: none"> <li>• Kostenlos*</li> <li>• Filtermöglichkeiten</li> <li>• Keine Begrenzung des Zeitintervalls</li> <li>• Verfolgung/Protokollierung von gelöschten Tweets</li> </ul>  | <ul style="list-style-type: none"> <li>• Kostenlos</li> <li>• Vielzahl an Abfrage- und Filtermethoden</li> <li>• Nachträgliche Datenabfrage</li> </ul>   | <ul style="list-style-type: none"> <li>• Umfassende Daten (hinsichtlich Zeitraum)</li> <li>• Usability</li> <li>• Minimaler Aufwand</li> <li>• Keine Programmierkenntnisse notwendig</li> <li>• Bereits strukturierte und angereicherte Daten</li> </ul> |
| <b>NACHTEILE</b>     | <ul style="list-style-type: none"> <li>• Bandbreiten-Limitierung, die vor allem bei Großereignissen schnell erreicht wird</li> <li>• Keine historischen Daten</li> <li>• Favorites, Retweets schwer zu erheben</li> <li>• Kein Anspruch auf Vollständigkeit</li> <li>• Keine Informationen über Gesamtvolumen für Einordnung</li> <li>• Prozess muss kontinuierlich laufen, jede Unterbrechung bedeutet Datenverlust</li> </ul> | <ul style="list-style-type: none"> <li>• Restriktive Datenausgabe</li> <li>• Kurzer Zeithorizont</li> <li>• Komplexität der Abfragestruktur</li> <li>• Kein Anspruch auf Vollständigkeit</li> <li>• Keine Informationen über Gesamtvolumen auf Twitter für Einordnung</li> </ul> | <ul style="list-style-type: none"> <li>• Kostenpflichtig und kostspielig</li> <li>• Vollständigkeit des Tweet-Sets nicht gewährleistet</li> <li>• Keine Angaben über Reliabilität</li> <li>• <i>Black Box</i> – Erhebungsverfahren unbekannt</li> </ul>  |

\*) bei Verwendung der standardmäßigen Streaming API (Spritzer)

## 4.2 Systeme der Datenverwaltung

Dieses Kapitel behandelt die Speicherung und Verwaltung der Daten, die durch die im vorigen Kapitel gezeigten verschiedenen Ansätze gewonnen wurden. Hierfür stehen sehr unterschiedliche Möglichkeiten zur Verfügung. Zuerst soll das Speichern in einfache Text-Dateien betrachtet werden, das bereits im vorangegangenen Kapitel kurz angesprochen wurde. Schließlich befasst sich Kapitel 4.2.2 mit Datenbanksystemen, die ein systematischeres und verlässlicheres Sichern der Daten ermöglichen. Hierbei werden zunächst zwei unterschiedliche Datenbank-Architekturen skizziert, um schließlich beispielhaft den Speicherprozess zum Datenbanksystem *MongoDB* zu zeigen. Abschließend folgt ein Vergleich aller betrachteten Alternativen.

### 4.2.1 Speicherung in Textdateien

Die wohl einfachste Methode zum Speichern von Tweets ist die Verwendung von einzelnen Dateien, die die Tweets deserialisiert im JSON-Format oder bereits formatiert, gefiltert und umstrukturiert beinhalten. Für den Export von Twitter-Daten eignen sich TXT- oder JSON-Dateien, die Tweets im ursprünglichen JSON-Format umfassen, oder CSV-Dateien, die bereits selektierte Werte beinhalten. Python bietet die Möglichkeit, Tweets sowohl im JSON-Format abzuspeichern, als auch umstrukturiert im CSV-Format. Während die benötigten Programmteile zum Lesen der Daten bereits in Python integriert sind, muss jedoch ein zusätzlicher Konverter zur Ausgabe strukturierter JSON-Daten installiert werden, um Datenströme nahezu direkt in Textdateien zu schreiben. Es müssen lediglich die im Binär-Code vorliegenden JSON-Daten (BSON) aus der Twitter-API in ein strukturiertes JSON-Format deserialisiert werden. Dies erfolgt in der Regel mit einem `jsonpickle.encode`-Befehl. Die Codierung erfolgt schrittweise je Datensatz.

Der Auszug in Listing 9 zeigt ein einfaches Verfahren zum Abspeichern von Tweets im JSON-Format. In diesem Beispiel sollten Tweets mit dem Begriff „Tatort“ nachträglich über die REST APIs gesammelt werden. Nach der Definition des Datei-Namens und Formates wurde im Skript das Vorgehen bei neuen Tweets vorgegeben. Jeder neue Tweet wird im JSON-Format als neue Zeile dem Dokument angefügt, sodass jede Zeile einem Tweet entspricht. Jeder Schreibvorgang erhöht den Zähler um 1. Bei Erreichen des definierten Maximalwerts stoppt die Schleife den Speicherprozess.

*Listing 9:* Speicherung von Tweets in eine Textdatei (Auszug aus Listing 8)

```
[...]
exportfile = "tatort.txt"
maxTweets = 1000000
[...]
with open(exportfile, "w") as file:
    while tweetCount < maxTweets:
        try:
            [...]
            for tweet in newTweets:
                file.write(jsonpickle.encode(tweet._json,
                    unpicklable=False) + "\n")
            tweetCount += len(newTweets)
        [...]
    [...]
```

Die Begrenzung der Tweetsammlung ist in mehreren Hinsichten sinnvoll: Sie dient indirekt zur Kontrolle von Umfang und zeitlicher Relevanz, sodass beispielsweise nur die letzten/neuesten 1.000 Tweets gesammelt werden. Letztlich steuert eine Begrenzung auch die Dateigröße, was sich wiederum auf die erforderliche Rechenleistung auswirkt. Größere Dateien mit mehreren Gigabyte erfordern für die Analyse und Verarbeitung der gesammelten Daten eine bessere Systemausstattung. Zudem gibt es je nach Dateisystem unterschiedliche Anforderungen hinsichtlich der maximal erlaubten Dateigröße: *FAT32*, ein gängiger Standard von USB-Sticks, erlaubt beispielsweise nur Dateigrößen bis etwa 4 Gigabyte. Sollte der Datenumfang sehr groß und nur schwer kalkulierbar sein, weil Tweets über einen längeren Zeitraum mit Hilfe der Streaming API gesammelt werden, besteht die Möglichkeit, den Datensatz auf mehrere Dateien zu verteilen. Dies soll im folgenden Beispiel näher erläutert werden.

#### 4.2.1.1 Anwendungsbeispiel: Speichern von Tweets in JSON- und CSV-Dateien

Fällt bei der Wahl der Speichermethode von Twitter-Daten die Entscheidung auf Textdateien, ist es in vielen Fällen (wie bereits oben erwähnt) sinnvoll, die Daten auf mehrere Dateien zu verteilen. Auch hierfür kann ein Zähler eingesetzt werden, um die Aufteilung der Daten zu steuern. Listing 10 (nächste Seite) zeigt den Programmauszug eines Skripts zum sequentiellen Speichern in einzelnen JSON-Dateien. Dazu wurde die Klasse `tweepylistener` aus Listing 5, die der Erhebung von Tweets mit dem Inhalt „obama“ diente, um einige Parameter erweitert. Die



Datei-Benennung erfolgt in diesem Anwendungsfall mit dem Zeitstempel der jeweiligen Dateierstellung, was ein Suchen beziehungsweise Eingrenzen von Tweets auf bestimmte Dateien anhand der Zeitangabe ermöglicht. Zusätzlich wird ein Datei-Präfix eingefügt, das in diesem Fall dem Suchterm entspricht. Jeder neue Tweet erhöht den Zählerstand counter. Nach einem vorher bestimmten, maximalen Zählerstand (hier 1000) legt das Skript alle weiteren Tweets in einer neuen Datei ab, bis wiederum das definierte Limit erreicht ist.

Neben dem direkten JSON-Export erlaubt Python auch eine Konvertierung in das CSV-Format. Hier besteht die Möglichkeit, bereits vorher relevante Datenfelder auszuwählen und diese dann je Tweet zeilenweise zu speichern. Dies ist von Vorteil, wenn zum einen keine Möglichkeiten zur nachträglichen Verarbeitung von JSON-Daten bestehen oder wenn zum anderen bereits vor Erhebung sicher ist, dass nur bestimmte Datenfelder benötigt werden. Eine Beschränkung auf wesentliche Daten verringert den Datenumfang deutlich, da einige Entities Informationen beinhalten, die für Auswertungen unerheblich sind, wie beispielsweise grafische Angaben zum Profil.

*Listing 10: Erweiterter Prozess zum Speichern in mehreren Textdateien*

```
import tweepy, time, sys, json
[...]
class tweepylistener(tweepy.StreamListener):
    def __init__(self, api = None, prefix = "obama"):
        self.api = api or API()
        self.counter = 0
        self.prefix = prefix
        # Definiere Zieldatei für Tweet-Daten
        self.output = open("<Datei-Pfad>" + prefix + "." +
                           time.strftime("%Y%m%d-%H%M%S") +
                           ".json", "w")
        # Definiere Zieldatei für Daten gelöschter Tweets
        self.deleted = open("geloescht.txt", "a")
    [...]

    # Bestimme Vorgehen in unterschiedlichen Situationen
    # Fall 1: neuer Status-Tweet
    def on_status(self, status):
        try:
            self.output.write(status + "\n")
            self.counter += 1
```

```

if self.counter >= 1000:
    self.output.close()
    self.output = open("<Datei-Pfad>" + self.prefix +
                       "." + time.strftime("%Y%m%d- \
                       %H%M%S") + ".json", "w")
    self.counter = 0
    sys.stdout.write(time.strftime("%Y%m%d-%H%M%S") +
                     ">> Neue Datei erstellt \n")
    return
[...]
```

Ein Export *aller* Felder im CSV-Format ist allerdings aufgrund des leistungs- und zeitintensiven Prozesses zur Extrahierung, Decodierung und Restrukturierung der Daten nicht sinnvoll. Zwar ist das CSV-Format mit seiner Tabellen-Struktur deutlich kompakter als das JSON-Format, bei dem jedes Daten-Feld einzeln definiert werden muss. Jedoch ist ein Zusammenfügen von Datensätzen im JSON-Format wesentlich komfortabler und sicherer. Bei CSV muss die Reihenfolge und Vollständigkeit der Felder unbedingt eingehalten werden, da sonst die Zuordnung der Werte nicht mehr stimmt. Fehlende Werte erhalten einen Leereintrag in einer Zelle, um die feste Struktur der Daten zu erhalten. Dagegen sind JSON-Daten unsortiert und können in ihrem Umfang unvollständig sein: Nur vorhandene Werte werden gespeichert, die Reihenfolge ist unerheblich. Dadurch können JSON-codierte Daten einfach aneinandergesetzt werden, während bei CSVs zuerst die Reihenfolge der Daten und Konsistenz der Struktur überprüft werden muss.

Das Skript in Listing 11 extrahiert zur Veranschaulichung einzelne Tweet-Entities aus Streaming-Daten schreibt diese in eine CSV-Datei. Hierfür wurde der Code aus Listing 10 abgeändert und erweitert. Zuerst wird eine CSV-Datei erstellt und die Spaltennamen definiert. Anschließend wird das Vorgehen bei einem neuen Tweet vorgegeben: Ein *Writer* soll die Daten im CSV-Format in die vorher angegebene Output-Datei schreiben und einzelne Werte mit einem Semikolon trennen. Zudem definiert der Parameter `dialect="excel"` eine Formatierung, die ein späteres Öffnen und Interpretieren durch Excel erlaubt. Mit dieser Option legt das Skript alle für die Dateistruktur relevanten Parameter automatisch fest. Schließlich erfolgt die Auswahl einzelner Entities, die aus dem JSON-Code extrahiert werden. Der Tweet-Inhalt wird zudem in das weit verbreitete Format UTF-8 kodiert, welches sich zur Darstellung sprachspezifischer Zeichen (z.B. chinesischer Zeichen oder Umlaute) eignet. Schließlich hängt das Skript die Daten zeilenweise an das CSV-Dokument an. Ein Zähler erzeugt nach jeweils 1.000 Dokument-Zeilen eine neue Datei.

Listing 11: Export gesammelter Tweets in Semikolon-getrennte CSV-Dateien

```
[...]
class tweepylistener(tweepy.StreamListener):

    def __init__(self, api = None, prefix = "apple"):
        [...]
        # Definiere Zieldatei für Tweet-Daten
        self.output = open("<Datei-Pfad>" + prefix + "." +
                           time.strftime("%Y%m%d-%H%M%S") +
                           ".csv", "w")
        # Schreibe Spaltennamen für CSV-Datei
        self.output.write("Tweet;Zeit;Txt;User;UID;Name"+"\\n")

    def on_data(self, data):
        [...]
        def on_status(self, status):
            writer = csv.writer(self.output, delimiter=";",
                                dialect="excel")
            # Definiere einzelne Tweet-Entities
            t_id = json.loads(status)["id"]
            t_tme = json.loads(status)["created_at"]
            t_txt = json.loads(status)["text"].encode("utf-8")
            t_uid = json.loads(status)["user"]["id"]
            t_usr = json.loads(status)["user"]["screen_name"]
            # Fülle Zeile mit den Werten
            writer.writerow([t_id, t_tme, t_txt, t_uid, t_usr])
            self.counter += 1
            if self.counter >= 1000:
                self.output.close()
                self.output = open("<Datei-Pfad>" + self.prefix + "." +
                                    time.strftime("%Y%m%d-%H%M%S") +
                                    ".csv", "w")
                self.counter = 0
                sys.stdout.write(time.strftime("%Y%m%d-%H%M%S") +
                                   ">> Neue Datei erstellt" + "\\n")
            return
        [...]
```

#### 4.2.1.2 Bewertung der Speicherung in Text-Dateien

Zusammenfassend ist der Speicheransatz über Textdateien ein geeignetes Mittel, um schnell und ohne größeren Aufwand eine kleine bis mittlere Datenmenge zu speichern. Da Python bereits eine JSON- und CSV-Funktionalität integriert hat, sind keine zusätzlichen Installationen notwendig. Der eigentliche Speicher-Prozess ist kurz und erlaubt eine direkte Einflussnahme auf die Verarbeitung der Daten bis zum Schreibbefehl: So können einzelne Datenfragmente extrahiert oder umstrukturiert werden. Die Methode ermöglicht zudem die Wahl eines gewünschten Endformats: Der Export in das gängige CSV- oder TXT-Format erlaubt hier nicht nur ein direktes Bearbeiten, Suchen, Extrahieren oder Einfügen von Datensätzen in der Datei, sondern auch eine schnelle Weitergabe von Daten an andere Personen. Zusätzlich kann die Dateigröße mit der Festlegung der Größe von Datensätzen kontrolliert werden. So können auch Dateien in kleinerem Umfang erzeugt werden, um sie zum Beispiel per E-Mail oder auf einem USB-Stick weiterzugeben. Die einfache und offene Dateistruktur mit einer oder mehreren Dateien im CSV-, JSON- oder TXT-Format erlaubt das direkte Lesen und Bearbeiten einzelner Datensätze. Aufgrund der klaren Struktur (jede Textzeile enthält alle Informationen zu einem Tweet) erhält man schnell einen guten Überblick über den Datenbestand. Letztlich ermöglichen fertige Textdateien auch einen schnellen Import in Statistik-Programme, wie *SPSS* oder *Stata*, die Informationen direkt aus CSV- oder JSON-Dateien auslesen können.

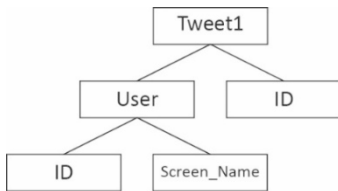
Jedoch eignet sich diese Methode nicht für Erhebungen mit sehr großem Datenumfang und/oder langer Dauer. Hier würden, je nach Wahl des `counter`-Parameters, wenige sehr große oder viele kleine Dateien entstehen. Zu große Dateien können unter Umständen nicht von allen Programmen eingelesen werden: So verarbeitet *Microsoft Excel* nur CSVs mit maximal etwa 65.000 Zeilen, *Microsoft Access* nur Daten bis zu 2 Gigabyte. Sehr große Text-Dateien mit mehreren Gigabyte an Informationen benötigen zudem leistungsfähige Rechner mit ausreichend großem Arbeitsspeicher (optimal sind 16 GB und mehr). Zu kleine Dateien führen mit zunehmender Dateizahl wiederum zu einer unübersichtlichen Dateilandschaft. Zudem stellt sich die Frage, wie die in den Dateien enthaltenen Informationen verknüpft und verarbeitet werden sollen: Einzelne Auswertungen je Datei oder aggregierte Auswertungen über alle Dateien? Für umfassende Analysen müssten alle Einzeldateien zunächst nachträglich zusammengefügt werden.

Ein großer Nachteil von Textdateien ist die mangelnde Zugriffsmöglichkeit auf die Daten. Während der Erhebung ist die momentan beschriebene Datei für Änderungen gesperrt, sodass beispielsweise nicht die ID des zuletzt hinzugefügten Tweets ausgelesen oder die Zuverlässigkeit des Programms überprüft werden

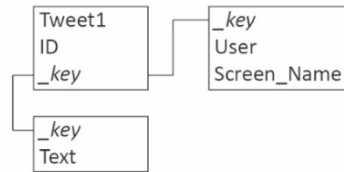
kann. Nach Abschluss des Schreibprozesses besteht keine Möglichkeit für einen simultanen Zugriff auf die Daten durch mehrere Personen. Zudem kann die zugreifende Person nur einen Prozess (Suchen, Kopieren, Überschreiben etc.) gleichzeitig durchführen. Funktionen einer komplexeren Datenverwaltung, wie ein Umstrukturieren, Aggregieren oder Filtern von Daten ist, wenn überhaupt, nur eingeschränkt möglich. Hierfür eignen sich besonders Datenbanksysteme, die das nächste Kapitel betrachtet.

### 4.2.2 Datenbank-Systeme

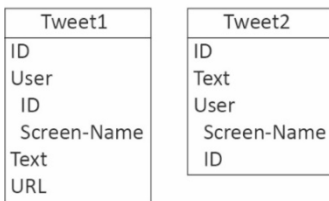
Für die Speicherung von Tweets stehen grundsätzlich alle Arten von Datenbanken zur Verfügung. Diese untergliedern sich in mehrere Datenbanksysteme, wie unter anderem in hierarchische, relationale und NoSQL-Datenbanken, und unterscheiden sich teils grundlegend in ihrer Datenstruktur (siehe Abbildung 10).



Hierarchisches DB-System



Relationales DB-System



Dokumentorientiertes NoSQL-DB-System

Abbildung 10: Vergleich mehrerer Datenbank-Systeme. Eigene Darstellung.

Bei hierarchischen DB-Systemen, der ältesten Datenbank-Architektur, sind alle Datensätze (*Records*) über eine Baumstruktur mit über- und untergeordneten Werten verbunden. Die Verknüpfung von Daten basiert auf Eltern-Kind-Beziehungen. Bei den sehr weit verbreiteten relationalen Datenbanksystemen (RDBMS), die häufig mit dem Begriff *SQL* zusammengefasst werden, sind Daten in der Regel über mehrere Tabellen verteilt und anhand von Schlüsseln miteinander verknüpft. So besteht ein *Record* unter Umständen aus mehreren Einträgen in verschiedenen Tabellen, wobei ein Eintrag üblicherweise aus einer Tabellenzeile (*Tupel*) mit mehreren Spalten (*Attributen*) besteht. Ein für jeden *Record* einmaliger Schlüssel stellt schließlich die Verbindung zwischen den einzelnen *Tupeln* dar. Vor der Dateneinspeisung in eine *SQL*-Datenbank muss zunächst jeder *Record* geparsed werden, um Daten in eine verwertbare Struktur und Formatierung zu überführen (Tugores & Colet, 2013, S. 21). Dieser Schritt wird bei einer größeren Datenmenge zeit- und rechenintensiv, weshalb andere Datenbank-Architekturen, wie *NoSQL*, unter Umständen eine bessere Eignung für das Sammeln großer Datenmengen aus dem Social Web haben.

NoSQL-Systeme („not only SQL“) erweitern die Funktionsweise von RDBMSs um mehrere Datenstruktur-Konzepte, auf denen die Daten verteilt und verknüpft sein können. Die wichtigsten Vertreter davon sind: Key-Value- und Graphen-Datenbanken sowie spaltenorientierte und dokumentorientierte Datenbanken. Key-Value-Datenbanken verwalten *Tupel*, die aus einem Schlüssel und dem dazu gehörenden Wert bestehen. Graphen-Datenbanken verknüpfen Datensätze über Knoten und Kanten. Spaltenorientierte Datenbanken (auch *Wide Column Stores*) haben, im Gegensatz zu RDBMS, eine dynamische Spaltenzahl. Es müssen also für einzelne Datensätze nicht alle Spalten definiert sein. Dokumentorientierte Datenbanken bündeln und speichern einzelne Datensätze in Dokumenten. Dokumente bestehen aus einer geordneten Liste von Key-Value-Paaren, wobei der Daten-Umfang dynamisch ist. (Trelle, 2014, S. 3)

Relationale Datenbanksysteme sind aufgrund ihrer Verbreitung und Funktionsfülle ideal für eine Vielzahl an Problemen. Bei der Wahl des Systems für die Sammlung großer Datensätze aus dem Internet sollten jedoch vor allem die jeweiligen Eigenschaften hinsichtlich Datenstruktur und Geschwindigkeit berücksichtigt werden. Durch das Web 2.0 gab es drastische Änderungen bei den Anforderungen an Datenbanksystemen: Sehr große Datenmengen sollen mit einem einfachen Schema bei geringer Datenkonsistenz und kurzen Laufzeiten abgespeichert werden (Trelle, 2014, S. 2). Zudem sollte die Möglichkeit bestehen, Speichersysteme problemlos zu erweitern (Skalierung).

Vergleicht man die oben genannten Architekturen, wird schnell die hervorragende Eignung von dokumentorientierten NoSQL-Systemen für die Sammlung von Twitter-Daten ersichtlich. NoSQL-Systeme erlauben eine horizontale Skalierung – im Lauf des Sammelprozesses können problemlos weitere Server zur Erweiterung des Speicherplatzes hinzugefügt werden (Chodorow, 2013, S. 4). Dies ist hilfreich, da eine Schätzung des erwarteten Datenvolumens (im Vergleich zu Befragungen) sehr schwer ist.

Im Vergleich zu relationalen Datenbanken mit einer Verteilung von Datenbündeln auf mehrere Tabellen bildet in dokumentorientierten Datenbanken ein Dokument genau einen Tweet-Record ab, der die von der Twitter-API übermittelten Key-Value-Paare und Arrays enthält. Somit entfallen Querverweise zu anderen Teilen eines Datensatzes und damit zusätzliche Datenbankabfragen. Dies beschleunigt den Abfrageprozess enorm und minimiert die Systemauslastung. Tugores und Colet (2013) verglichen die Schreibleistung des NoSQL-Systems MongoDB mit MySQL und erkannten bei ersterem eine deutlich höhere Verarbeitungsgeschwindigkeit bei geringerem Leistungsbedarf.

Die Datenstruktur von NoSQL-Systemen ist zudem eine abgewandelte Version des JSON-Formates, welches Twitter bei der Datenübermittlung nutzt. Dokumentorientierte Datenbankssysteme haben eine geringe Konsistenz-Anforderung an Daten: Datensätze können in ihrem Umfang und Aufbau variieren. Demgegenüber müssen bei SQL-Systemen Daten-Records einem vorher definierten Schema entsprechen. Verschachtelte Daten aus Arrays müssen in separate Tabellen geschrieben werden, die wiederum mit einem Schlüssel verknüpft sind. Während bei Document Stores einzelne Records schnell und einfach identifiziert und exportiert werden können, müssen diese bei SQL-Datenbanken zuerst anhand der Verknüpfungen zusammengefügt werden.

Aufgrund der guten Eignung von No-SQL-Systemen für die Speicherung großer Mengen von Internet-Daten verwendet diese Arbeit das dokumentorientierte Datenbank-System *MongoDB*, welches das folgende Kapitel kurz vorstellt.

#### 4.2.2.1 MongoDB

MongoDB wurde 2007 entwickelt und wird seitdem als Open Source System angeboten. Das Datenbank-System besteht in seiner einfachsten Form aus einem Server (Host), der die Daten speichert und verwaltet, sowie einem oder mehreren Clients, die Daten in die Datenbank einspeisen oder abrufen. Das System ist skalierbar, sodass zu einem späteren Zeitpunkt weitere Hosts und Clients hinzugefügt

werden können. Eine Benutzerverwaltung steuert den Zugriff und die Berechtigungen einzelner Nutzer.

MongoDB strukturiert Daten anhand von Objekten, die in Collections gesammelt werden. Eine Collection ist eine Sammlung von Objekten, bei der, im Gegensatz zu SQL-Tabellen, eine Schemafreiheit besteht. Es wird also kein festes Schema mit Namen, Typen und Reihenfolge für einzelne Felder vorgegeben, sondern Daten relativ unstrukturiert gespeichert. Mehrere Collections werden wiederum in einer Datenbank gebündelt. Der Schreibzugriff auf Collections beziehungsweise Datenbanken ist dabei nicht auf einen Nutzer beschränkt, wie das bei Text-Dateien der Fall ist. Dadurch ist es möglich, dass mehrere Prozesse parallel Twitter-Daten sammeln und in verschiedene Collections speichern. Technisch besteht zudem die Möglichkeit, in einem Abfrage-Prozess über eine Weiche spezifische Tweets, z.B. anhand eines Schlagwortes, auf bestimmte Collections zu verteilen. MongoDB beinhaltet außerdem mehrere Funktionen zum Filtern, Sortieren, Umstrukturieren, Importieren und Exportieren sowie grundlegende Möglichkeiten zum Analysieren von Daten.

Vor dem Speichern in MongoDB muss die Datenbank zuerst installiert und eingerichtet werden. Eine Trennung von Datenbank-Server und Client ist dabei sinnvoll<sup>22</sup>. Des Weiteren empfiehlt sich das Anlegen mehrerer Nutzer, die entweder Zugriff auf die Verwaltungsoptionen oder Lese-/Schreibzugriff auf bestimmte Datenbanken haben. Das gewährleistet nicht nur eine allgemeine Sicherheit der Daten vor Fremdzugriff, sondern verhindert zusätzlich, dass Prozesse versehentlich (bei falscher Angaben des Collection-Parameters) Daten in eine falsche Collection schreiben. Zudem ist es ratsam, das Datenbankprogramm als Dienst einzurichten, der bei Bedarf automatisch den Server-Prozess startet.

Sobald der Datenbank-Server eingerichtet ist, kann der Zugriff auf MongoDB von jedem PC erfolgen, der einen MongoDB-Client installiert hat. Grundsätzlich besteht die Möglichkeit, per Kommandozeile die Datenbank zu verwalten, Daten zu lesen, zu sortieren und filtern. Zusätzlich gibt es eine Vielzahl leicht bedienbarer Management-Systeme mit grafischen Benutzeroberflächen. Einige sind dabei bei eingeschränktem Funktionsumfang kostenlos verfügbar, wie beispielsweise

---

<sup>22</sup> Für eine Gewährleistung von Datensicherheit, Datenschutz und Performance empfiehlt sich immer eine Trennung von Datenbank und Nutzer.



*Mongo Management Studio*<sup>23</sup> (siehe Abbildung 11), *MongoVUE*<sup>24</sup> oder *3T MongoChef*<sup>25</sup>. Über diese Systeme können gängige Prozesse, wie der Import oder Export, sowie die Suche oder das Filtern von Daten einfach ausgeführt werden.

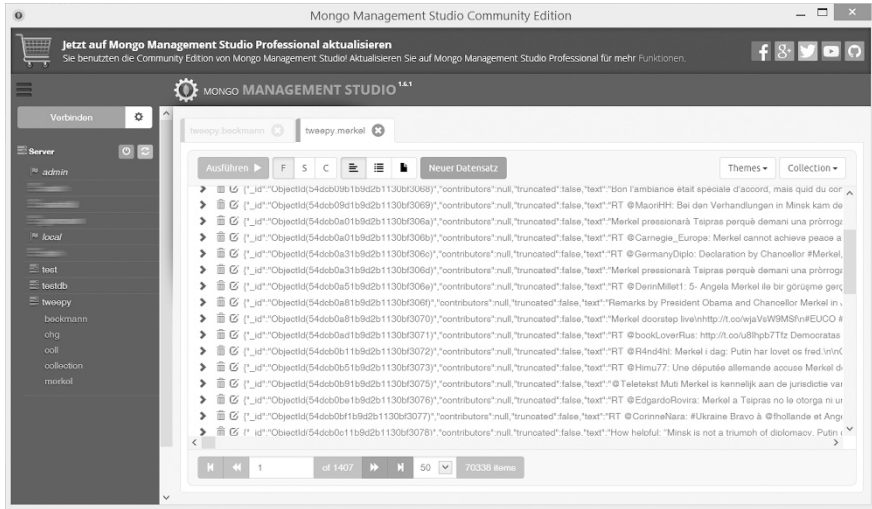


Abbildung 11: Benutzeroberfläche des Mongo Management Studio.

Für eine Verbindung per Kommandozeile muss zuerst im Shell-Fenster der Installationsordner von MongoDB aufgerufen werden, beziehungsweise dessen Unterverzeichnis `/bin`, in dem sich die Startdateien befinden. Zum Start werden weitere Parameter, wie die IP-Adresse des Datenbank-Servers, der Port und Zugangsdaten angegeben. Mit Hilfe der Befehle `show dbs` und `show collections` können Datenbanken und – nach Selektion per `use`-Befehl – die jeweiligen Collections angezeigt werden.

Der Code aus Listing 12 zeigt das übliche Vorgehen zum Aufrufen einer Collection. Bei Verbindung mit der Datenbank findet eine Autorisierung (hier mit einem Administrator-Konto) statt. Nach einer Auflistung der Datenbanken und der Auswahl der Datenbank `tweepy`, werden deren Collections abgerufen. In diesem Fall soll zudem die Twitter-ID des zuletzt gespeicherten Tweets angezeigt werden.

<sup>23</sup> <http://www.litixsoft.de/mms/>

<sup>24</sup> <http://www.mongovue.com/>

<sup>25</sup> <http://3t.io/mongochef/>

Dabei werden mehrere Abfrage-Operatoren miteinander kombiniert: `find`, `limit` und `sort`. Die Tweet-ID ist im Schlüssel `id_str` hinterlegt. Da MongoDB jedem Datenbank-Objekt nochmals eine eigene ID zuweist (Feld „`_id`“), soll diese bei der Ausgabe ausgeblendet werden, um Verwechslungen zu vermeiden.

Die zweite geschweifte Klammer des Operators `find` definiert die Ausgabe der Felder. Um nur einen Datensatz zu erhalten, wird ein Limit von 1 vorgegeben. Zudem erfolgt eine absteigende Sortierung anhand der Objekt-ID, um den neuesten Tweet (mit der höchsten ID) zu erhalten. Listing 12 stellt nochmals alle vorgestellten Befehle und deren Ergebnisse dar.

### Listing 12: Einfache Datenabfrage bei MongoDB

```
C:\<Pfad>\bin>mongo -host <IP-Adresse>
MongoDB shell version: 2.6.7
connecting to: <IP-Adresse>:<Port>/mongo
> use admin
switched to db admin
> db.auth("<DB-Benutzer>", "<DB-Passwort>")
1
> show dbs
admin          0.078GB
local          0.078GB
tweepy        0.953GB
> use tweepy
switched to db tweepy
> show collections
chg
merkel
system.indexes
> db.merkel.find({}, {id_str: 1, _id:
0}).limit(1).sort({$id:-1})
{ "id_str" : "565872318040514560" }
>
```

Diese Methode ist die schnellste Möglichkeit zur Abfrage spezifischer Daten bei MongoDB und sollte, wenn möglich, immer auf diese Weise angewendet werden. Ein Sortieren der Daten über Verwaltungssoftware wie Mongo Management Studio benötigt bei großem Datenumfang sehr viel Arbeitsspeicher. Ist nicht genügend Arbeitsspeicher vorhanden, bricht die Operation mit einer Fehlermeldung ab. Ein Grund für den erhöhten Speicher-Bedarf ist die Art des Such-Prozesses: Die

grafische Benutzeroberfläche muss die sortierten Tweets auch gleich am Bildschirm anzeigen, weshalb immer die kompletten Objektdaten (also alle Meta-Daten) verarbeitet werden müssen. Der Befehl über das Shell-Fenster beschränkt die Suche bereits vorher auf das Feld `id_str`. Jedoch besteht auch bei vielen Verwaltungssystemen die Möglichkeit, Such-, Filter- und Sortierbefehle in ein Abfragefeld einzugeben, sodass nicht zwingend ein Zugriff per Shell erforderlich ist.

MongoDB erlaubt mit ähnlich strukturierten Befehlen auch die Selektion, Aggregation beziehungsweise Restrukturierung von Daten. Diese Funktionen betrachtet Kapitel 4.3.2 genauer. Das folgende Beispiel skizziert den Ablauf eines Speicherprozesses.

#### 4.2.2.2 Anwendungsbeispiel: Speichern von Tweets in MongoDB

Nachdem bereits eine Daten-Abfrage auf dem MongoDB-Server getätigt wurde, zeigt dieses Kapitel die Abwandlung des bestehenden Skriptes aus Listing 10 zum Speichern von Tweets. Für Python gibt es ein Programmpaket namens *pymongo*, das den Datenverkehr über die MongoDB Schnittstelle verwaltet. Listing 13 stellt die übliche Konfiguration der Datenbankverbindung dar. Zuerst werden die IP-Adresse des Datenbank-Servers sowie, wenn eingerichtet, Benutzername und Passwort angegeben. Zusätzlich folgt der Name der gewünschten Datenbank. Bei Verbindungsproblemen oder fehlerhaften Anmeldedaten erscheint eine entsprechende Fehlermeldung und das Skript beendet die Verbindung zu MongoDB.

Schließlich wird das Vorgehen bei einem neuen Status definiert. In diesem Fall schreibt das Skript jeden Tweet, der über die Streaming API übermittelt wird, direkt in die Datenbank. Möglich wäre aber auch eine vorherige Definition des Datenumfangs, sodass der Insert-Befehl nur bestimmte Werte (wie Name, Tweet-ID und Datum) umfasst. Da außerdem die Möglichkeit besteht, verschiedene Werte in unterschiedliche Collections zu schreiben, ist bei jedem Insert die Angabe einer Ziel-Collection erforderlich. Der Parameter `safe=True` stellt sicher, dass die Datenbank bei jedem Schreibprozess auf eine Reaktion, genauer gesagt auf die Bestätigung einer korrekten Ausführung des Speicherprozesses, wartet. Dies verhindert einen Datenverlust durch unvollständige Schreibbefehle (bei Verbindungsproblemen oder Überlastung). Bei Fehlern gibt das Skript eine Fehlermeldung aus und unterbricht die Verbindung.

Listing 13: Auszug eines Prozesses zum Speichern von Tweets in MongoDB

```

import pymongo
[...]
class tweepylistener(tweepy.StreamListener):
    def __init__(self, api = None):
        self.api = api or API()
        super(tweepy.StreamListener, self).__init__()

    try:
        connection = pymongo.MongoClient("<IP-Adresse>")
        connection.tweepy.authenticate("<DB-Benutzer>",
                                       "<Passwort>")

        print "Verbindung zu MongoDB erfolgreich"
        self.db = connection.<DB-Name>

    except ConnectionFailure, e:
        sys.stderr.write("keine Verbindung möglich: %s" % e)
        sys.exit(1)connection = pymongo.MongoClient("<IP- \
Adresse>")

[...]
def on_status(self, status):
    try:
        self.db.<Collection-Name>.insert(json.loads(status),
                                         safe=True)

```

#### 4.2.3 Vergleich der Systeme zur Datenverwaltung

Beide Arten der Datensicherung haben Vor- und Nachteile. Die Speicherung in Einzeldateien ist einfach auszuführen und bedarf keiner speziellen System-Konfiguration. Die erzeugten Dateien sind direkt einlesbar und können ohne Umwandlung weiterverarbeitet oder weitergegeben werden. Es bedarf keiner Installation weiterer Software. Jedoch wird die Dateistruktur gerade beim Sammeln großer Datenmengen schnell unübersichtlich, da große Datenmengen systembedingt auf mehrere Dateien verteilt werden müssen. Für eine vollständige Analyse müssten diese Datenfragmente vorher wieder zusammengefasst werden. Des Weiteren besteht ein hohes Risiko des Datenverlustes: Es sind regelmäßige Backups der Dateien nötig, welche aber nur bei bereits vollständig beschriebenen Dateien möglich sind. Sollte ein Systemfehler beim Schreibprozess auftreten, ist die momentan beschriebene Datei unter Umständen verloren. Somit eignet sich das Speichern in

Einzeldateien eher für Ad-hoc-Analysen oder zum Testen der gewählten Abfrageparameter der APIs.

Datenbanksysteme haben dagegen spezielle Anforderungen an Hard- und Software. Das System muss zuerst installiert und konfiguriert werden, indem beispielsweise Hosts und Clients, Benutzerkonten und Rechte eingerichtet werden. Dafür weist das System einen hohen Grad an Sicherheit hinsichtlich Zugriffsschutz und Datenbeständigkeit auf. Der Zugriff kann mithilfe von Benutzerrollen auf spezifische Datenbanken eingeschränkt und der Datenbestand durch Replikation der Daten auf mehrere Server gewährleistet werden. Zudem gestatten professionelle Datenbanksysteme simultane Lese- und Schreibzugriffe, sodass gleichzeitig sowohl mehrere Sammelprozesse, als auch Ad-hoc-Analysen des bestehenden Datenmaterials möglich sind. Das Speichern in Datenbanken eignet sich aufgrund des erhöhten Einrichtungsaufwands eher für eine größere, langfristig durchgeführte Datensammlung. Tabelle 6 fasst die Vor- und Nachteile der beiden Alternativen nochmals zusammen.

Tabelle 6: Vergleich der Möglichkeiten zur Datenspeicherung

|                      | TEXT-DATEIEN   | DOKUMENTORIENTIERTE DATENBANKEN  |
|----------------------|--|--|
| <b>EIGENSCHAFTEN</b> | <ul style="list-style-type: none"> <li>• JSON-, TXT-, CSV-Format</li> <li>• Record entspricht Textzeile</li> </ul>   | <ul style="list-style-type: none"> <li>• BSON/JSON-codiert</li> <li>• Record entspricht Dokument</li> </ul>  |
| <b>VORTEILE</b>      | <ul style="list-style-type: none"> <li>• Einfach &amp; schnell anwendbar</li> <li>• Mit „Bordmitteln“ des Betriebssystems realisierbar</li> <li>• Daten direkt einsehbar</li> <li>• Direkte Datenweitergabe möglich</li> </ul> | <ul style="list-style-type: none"> <li>• Datensicherheit und Zugriffsschutz</li> <li>• Gute Performance</li> <li>• Multi-User, Multi-Thread</li> <li>• Grundlegende Analyse- und Aggregations-Funktionen bereits integriert</li> </ul> |
| <b>NACHTEILE</b>     | <ul style="list-style-type: none"> <li>• Mangelnder Schutz vor Datenverlust</li> <li>• Nur ein Lese- oder Schreibprozess gleichzeitig</li> <li>• Bei großem Datenumfang unübersichtliche Einzeldateistruktur</li> </ul>        | <ul style="list-style-type: none"> <li>• Installation und Konfiguration zwingend notwendig</li> <li>• Unter Umständen große Hardware-Anforderungen</li> <li>• Programmierkenntnisse für Shell-Kommandos erforderlich</li> </ul>        |

Die Wahl des Systems wird, wie bereits bei der Wahl der Sammelmethode, vom Zweck beziehungsweise der Absicht der Datenverwendung beeinflusst. Generell sind immer Datenbanken zu empfehlen, da diese hinsichtlich Sicherheit, Beständigkeit, Leistung und Verwaltung einfachen Textdateien überlegen sind. Für simple Zwecke genügen jedoch auch Einzeldateien. Zudem besteht bei allen gängigen Datenbanksystemen die Möglichkeit, manuell gesammelte Daten (beispielsweise im CSV-, XML- oder XLS-Format) nachträglich zu importieren.

Ein Vorteil von MongoDB liegt in der Möglichkeit, eine Strukturierung und grundlegende Analyse der Daten bereits mit systemeigenen Funktionen durchführen zu können. Auf diese Analysemöglichkeiten und andere Python-basierte Auswertungen geht das folgende Kapitel genauer ein.

### 4.3 Methoden der Datenanalyse

Nachdem bereits Ansätze zum Sammeln und Speichern von Twitter-Daten besprochen wurden, folgt nun eine Betrachtung der Möglichkeiten zur Analyse dieser Daten. Dieses Kapitel stellt dabei keinen Anspruch an Vollständigkeit, sondern dient zur Darstellung und Bewertung mehrerer, gängiger Analyseansätze und Methoden. Zuerst wird auf etwaige Probleme bei der Datenanalyse und deren Ursachen eingegangen sowie Ansätze zur Vermeidung und Behebung thematisiert. Da in dieser Arbeit auf dem Datenbanksystem MongoDB basiert, zeigt Kapitel 4.3.2 zunächst die systemintegrierten Funktionen zur Datenverarbeitung und -analyse. Kapitel 4.3.3 befasst sich schließlich mit ausgewählten Methoden des *Natural Language Processing* (NLP) – der automatisierten Inhaltsanalyse der Computerlinguistik. Hierbei gilt zu beachten, dass dies nur eine kleine Auswahl gängiger Analyseverfahren ist, die spezifisch für das folgende Beispiel getroffen wurde.

Der Vergleich der Methoden erfolgt anhand eines Beispiels – der Analyse von Tweets über den Franken-Tatort<sup>26</sup>. Das Datenset umfasst dabei 37.556 vollständige Tweet-Records (inklusive aller Meta-Daten, siehe Abbildung 6 weiter vorne), die ex post mit der Search API (Suchterm: „tatort“) ermittelt wurden. Die Collection `tatort_tweets` in der Datenbank `tatort` beinhaltet in ihrer ursprünglichen Form Tweets unterschiedlichster Sprachen und dadurch mit großer Wahrscheinlichkeit auch irrelevante Tweets. Dies entspricht der üblichen Ausgangssituation bei Twitter-Analysen. Folglich empfiehlt sich in vielen Fällen ein vorgelagertes Bereinigen der Daten.

---

<sup>26</sup> „Der Himmel ist ein Platz auf Erden“, ARD, 12.04.2015, 20:15-21:45. Erfassung aller Tweets zwischen 05.04., 14:52, und 13.04., 15:34 Uhr.

### 4.3.1 Vorverarbeitung der Daten

Nach der Datensammlung und vor der eigentlichen Analyse der gesammelten Daten müssen diese zuerst aufbereitet werden, indem sie beispielsweise normalisiert, sortiert, aggregiert, gefiltert und korrigiert werden (Abbildung 12 zeigt den typischen Ablauf der Datenanalyse). Besonders bei großen Datensets sind mehrere vorgelagerte Filter zur Verringerung des Datenumfangs sinnvoll. So können beispielsweise irrelevante Werte der Tweet-Records gelöscht werden. Besonders Angaben über die grafische Gestaltung des Urheber-Accounts (Hintergrundbild, Farbe, etc.) sind für die meisten Auswertungen unbedeutend. Ebenso empfiehlt sich vor der Analyse das Aussondern anderer Sprachen oder Tweets außerhalb eines bestimmten Zeitraums. Die Filterung nach Sprache benötigt allerdings eine verlässliche Spracherkennung. Der Twitter-eigene Spracherkennungs-Dienst ist die momentan zuverlässigste Identifikations-Methode (Carter, Weerkamp & Tsagkias, 2013). Sofern jedoch Daten vor Einführung des Spracherkennungs-Algorithmus durch Twitter im März 2013 (Roomann-Kurrik, 2013) analysiert werden, bedarf es anderer Programme.

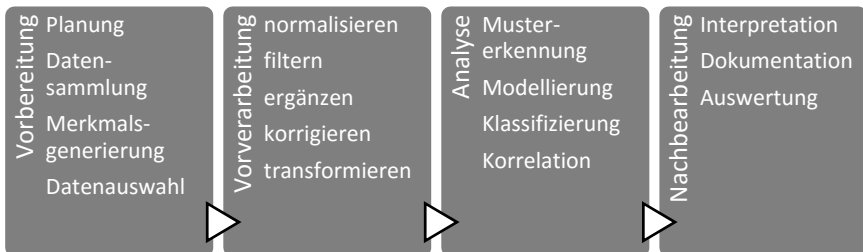


Abbildung 12: Prozess der Datenanalyse. In Anlehnung an Runkler (2000, S. 2).

Neben dem Twitter-eigenen Algorithmus gibt es auch andere Dienste, wie Googles kostenpflichtige *Translate API* (Google Inc., 2015), oder Skripte, wie *Langid* für Python (Lui & Baldwin, 2012). Alle verfügbaren Werkzeuge können keine hundertprozentig (oder annähernd) verlässliche Spracherkennung gewährleisten (Lui & Baldwin, 2014). Die Merkmale der Twitter-Sprache (Dialekt, Neologismen usw.), besonders die Vermischung mehrerer Sprachen in einem Tweet, aber auch die sehr limitierte Textlänge und Häufung von Abkürzungen stellen hier ein großes Hindernis dar (Carter et al., 2013). Da diese Eigenheiten nicht beein-

flussbar sind, besteht in der Zukunft nur die Möglichkeit, intelligentere und zuverlässigere (lernfähige) Algorithmen zu entwickeln. Momentane Lösungen bieten allerdings schon relativ verlässliche Ergebnisse.

Ein großes Problem, nicht nur bei der Inhaltsanalyse, könnten Spam-Tweets darstellen (McCord & Chuah, 2011). Das Veröffentlichen vieler (identischer) Tweets durch mehrere Accounts in kurzer Zeit kann beispielsweise dazu genutzt werden, eigene Hashtags (etwa von Aktionen) in die Trending Topics zu bekommen, um dadurch einfach eine größere Zielgruppe zu erreichen (Benevenuto, Magno, Rodrigues, & Almeida, 2010). Oftmals enthalten Spam-Tweets auch momentan populäre – und themenfremde – Hashtags und einen Link zu einer Webseite. Die Erkennung von Spam kann unter Umständen die Qualität des Datensatzes entscheidend erhöhen. Laut einem Bericht von Twitter, Inc. an die New Yorker Börsenaufsicht SEC schätzt das Unternehmen den Anteil an Spam-Bots, also Accounts, die automatisiert Spam verschicken, auf etwa 5 Prozent der monatlich aktiven Accounts (United States Securities and Exchange Commission [SEC], 2014).

Wenn neben reinen Häufigkeitsauszählungen auch quantitative Inhaltsanalysen durchgeführt werden sollen, muss das Tweet-Datenset noch detaillierter vorverarbeitet werden. Für eine Betrachtung des reinen Twitter-Textes, beispielsweise zur Ermittlung der häufigsten Begriffe eines Datensets, erscheint das Filtern von Mentions, Hyperlinks und Retweets als sinnvoll. Zudem sollten hier Satzzeichen und andere Sonderzeichen entfernt werden. Sofern eine Analyse nicht auf Satz-, sondern auf Wortebene erfolgt, empfiehlt sich der Ausschluss häufiger Begriffe, wie Artikel, Konjunktionen oder Präpositionen aus dem Datensatz. Man spricht hier von *Stoppwörtern*.

Eine Konvertierung in Kleinbuchstaben vereinfacht nicht nur eine spätere Lemmatisierung oder ein Stemming, sondern auch das Erkennen/Filtern von Begriffen. So müssen nicht alle Kombinationen möglicher Groß- und Kleinschreibungen in Begriffskatalogen gespeichert werden. Des Weiteren ist eine Umwandlung von Abkürzungen sinnvoll, da diese in der Internetsprache häufig vorkommen – besonders aufgrund der Beschränkung auf 140 Zeichen pro Tweet erscheint dies bei Twitter als wahrscheinlich. Für diese speziellen Bedarfssituationen muss die Umwandlung allerdings anhand manuell erstellter Lexika erfolgen. Im Gegensatz zu gängigeren Abkürzungen, wie „lol“ oder „ftw“, stellt die Konvertierung unkonventioneller Abkürzungen eine größere Herausforderung dar. Nutzerspezifische Sprache kann selbst durch ein Einlesen in themenspezifische Tweets nicht zwangsläufig erkannt werden. Sehr spezifische Abkürzungen oder falsch geschriebene Begriffe müssten manuell ersetzt und korrigiert werden. Zudem empfiehlt











sich eine Filterung und Korrektur nicht rein alphabetischer Begriffe (z.B. „w00t“, „n8“).

Neben Abkürzungen enthalten Tweets auch häufig Symbole, wie Emoticons, die bestimmte Emotionen ausdrücken. Zwar gibt es einen durch das Unicode Konsortium standardisierten Katalog von *Emojis* (Unicode, Inc., 2015), jedoch weisen Park, Barash, Fink und Cha (2013) auf kulturelle und persönliche Unterschiede in der Interpretation. Da Emoticons per Definition eine bedeutende Informationsquelle für vermittelte Emotionen sind, sollten diese bei der semantischen Analyse berücksichtigt werden. Um Verzerrungen durch anders interpretierte Smileys zu reduzieren, ist womöglich eine Beschränkung der Analyse auf die gängigsten Symbole – wie beispielsweise :-), :-(-;-) :-P :D – am sinnvollsten. Dennoch kann es auch hier zu unterschiedlichen Interpretationen kommen, da beispielweise jeder Hersteller von Betriebssystemen für Smartphones eigene Grafiken verwendet (siehe Tabelle 7).

Hinzu kommt die Tatsache, dass Twitter mittlerweile ein eigenes Set an *Twemojis* zur Verfügung stellt, das die von den jeweils unterschiedlichen Betriebssystemen verwendeten Icons in der Darstellung auf der Twitter-Webseite vereinheitlicht (Davidson, 2014). Das Beispiel unten zeigt zudem, dass die Gefahr besteht, die ursprüngliche Aussage des Emojis (hier: Schläfrigkeit) falsch zu deuten. Eine mögliche Interpretation wäre in diesem Fall auch ein trotziges Spucken oder Weinen.

Vorgelagerte Filter und Verarbeitungsschritte ermöglichen eine Reduktion des Datenumfangs und eine Verbesserung der Lesbarkeit von Texten. Zwar gehen dadurch möglicherweise Informationen verloren, dennoch kann die Aussagekraft des Inhalts durch Ausschluss irrelevanter Zeichen oder ganzer Text-Fragmente steigen. Zudem sinkt das Datenvolumen, was die Analyse beschleunigen kann. Vor der eigentlichen Analyse sollte deshalb immer eine Bereinigung erfolgen.

Tabelle 7: Unterschiede in der Darstellung von Emojis. Bildquelle: Unicode, Inc. (2015).

| HERSTELLER        | APPLE   | GOOGLE  | MICROSOFT   | TWITTER   |
|-------------------|---|---|---|---|
| SYSTEM            | iOS und OS X  | Android   | Windows   | Twemoji   |
| SCHLÄFRIGES EMOJI |  |  |  |  |
| WEINENDES EMOJI   |  |  |  |  |

#### 4.3.2 Verarbeitung und Analyse mit MongoDB

MongoDB stellt für das sogenannte *Preprocessing* eine Vielzahl an Möglichkeiten zur Manipulation und Analyse von Daten zur Verfügung. Diese eignen sich jedoch eher zur allgemeinen Strukturierung der Daten, als für die Textmanipulation oder detailliertere Analysen. So können Datensätze gefiltert oder gruppiert werden, während eine inhaltliche Verarbeitung von Texten (z.B. Stoppwort-Filter, Tokenisierung) nur auf Umwegen über andere Programme oder Skripte möglich ist.

Ein wichtiges Toolkit für die Verarbeitung auf Datensatz-Ebene sind die Methoden zur Datenaggregation, welche drei wesentliche Funktionen beinhalten: Abfragemethoden zur Aggregation, das *Aggregation Framework* sowie das *MapReduce*-Verfahren. Die drei Methoden unterscheiden sich im Allgemeinen vor allem in ihrer Performance und ihrem Funktionsumfang. Während die Abfragemethoden einen geringen Funktionsumfang, aber eine hohe Leistung hinsichtlich Verarbeitungsgeschwindigkeit vorweisen, bietet MapReduce eine sehr hohe Funktionalität (auch für Analysen) mit jedoch eingeschränkter Performance. Einen Mittelweg zwischen Leistung und Funktionalität geht das Aggregation Framework.

### 4.3.2.1 Abfragemethoden zur Aggregation

MongoDB stellt mit mehreren Abfragemethoden eine Grundfunktionalität zur Aggregation von Daten zur Verfügung. So können beispielsweise Häufigkeitsauszählungen von (gruppierten) Werten bereits mit den internen Funktionen der Datenbank durchgeführt werden. Da die Möglichkeiten zur Einschränkung der Suche über die REST APIs begrenzt sind, wurden im eingangs erwähnten Erhebungsfall auch Tweets mit dem Begriff Tatort erfasst, die deutlich vor dem Sendetermin lagen. Das Datenset besteht folglich aus Tweets zwischen dem 05.04. und 13.04. Für die Datenanalyse sollen zunächst alle Einträge vor dem 11.04. gefiltert werden. Um den Datenumfang anhand des Datums auf relevante Tweets zu begrenzen, muss das Feld `created_at` zunächst umformatiert werden. Twitter verwendet das Unicode-Format, welches aber nicht direkt kompatibel mit den Filteroperatoren von MongoDB ist. Listing 14 beinhaltet ein Skript zur Konvertierung der Zeitstempel in das `Datetime`-Format nach dem Standard ISO 8601 (International Organization for Standardization [ISO], 2015).

*Listing 14:* Konvertierung des Datumsformats

```
import pymongo, datetime
from pymongo import MongoClient

connection = pymongo.MongoClient("<DB-IP>")
connection.admin.authenticate("<DB-Benutzer>", "<Passwort>")
db = connection.<DB-Name>
collection = db.<Collection-Name>
tweets = collection.find({})
for tweet in tweets:
    tweetdate = tweet[u'created_at']
    if(type( tweetdate ) == unicode):
        newdate = datetime.datetime.strptime(tweetdate,
            '%a %b %d %H:%M:%S +0000 %Y')
        pointer = tweet[u'_id']
        collection.update({'_id': pointer},
            {'$set': {'created_at': newdate}})
    else:
        skip=1
    pass
print("Alle Daten konvertiert")
if skip==1:
    print("Mindestens 1 Datum wurde übersprungen")
```

Nach der Umwandlung des Datums besteht nun eine Filter-Möglichkeit nach Datum. Dadurch ergibt sich ein kleinerer Datensatz mit 20.310 Einträgen. Da der Tatort vornehmlich ein deutsches beziehungsweise deutschsprachiges Phänomen ist, interessieren vor allem Tweets auf Deutsch. Eine Sichtung des Materials ergab, dass manche Twitter-Nutzer momentan populäre Hashtags ohne thematischen Bezug in ihre Nachrichten einbauen, um vermutlich eine größere Beachtung zu finden (Spam). Dies konnte man vor allem bei fremdsprachigen Tweets erkennen. Auch deshalb werden für diese Analyse nur deutschsprachige Tweets analysiert. Twitter stellt für die Textsprache das Entity `iso_language_code` zur Verfügung. Twitters eigener Spracherkennungs-Algorithmus erkennt anhand mehrerer Parameter automatisiert die Textsprache und trägt diese als Code in das Feld ein. Aktuelle Studien bescheinigen eine relativ hohe Zuverlässigkeit, wenngleich unbereinigte Tweets eine nicht zu vernachlässigende Fehlerquote bei der Erkennung aufweisen (vgl. Kapitel 4.3.1).

#### *Listing 15:* Abfragemethoden zur Aggregation in der MongoDB-Shell

```
> db.tatort_tweets.count()
37556
> db.tatort_tweets.count({
...   created_at: {
...     $gte: ISODate("2015-04-11T12:00:00.000Z"),
...     $lt: ISODate("2015-04-13T12:00:00.000Z")}
... })
20310
> db.tatort_tweets.find({
...   created_at: {
...     $gte: ISODate("2015-04-11T12:00:00.000Z"),
...     $lt: ISODate("2015-04-13T12:00:00.000Z")},
...   "metadata.iso_language_code" : "de"
...   }).count()
18517
> db.tatort_tweets.group({
...   key: {"metadata.iso_language_code" : 1},
...   cond: {
...     created_at: {
...       $gte: ISODate("2015-04-11T12:00:00.000Z"),
...       $lt: ISODate("2015-04-13T12:00:00.000Z")}},
...   initial: {count:0},
...   reduce: function(curr,result){
...     result.count++;}
... })
```

```
[
  {
    "metadata.iso_language_code" : "de",
    "count" : 18517
  },
  {
    "metadata.iso_language_code" : "en",
    "count" : 739
  },
  [...]
]
```

Listing 15 beginnt mit einer einfachen Auszählung aller Dokumente in einer Collection. Diese allgemeine Abfrage liefert einen erstmaligen Überblick der Tweet-Anzahl. Hierfür wurde der Count-Befehl verwendet und schließlich um einen Datumsfilter erweitert. Um Tweets des relevanten Zeitraums in deutscher Sprache zu erhalten, wird Count mit dem Filter Find ergänzt, der die Abfrage des ursprünglichen Befehls ersetzt. Der Datensatz beinhaltet demnach 18.517 Tweets in deutscher Sprache. Zudem besteht die Möglichkeit, Daten anhand eines Schlüssels zu gruppieren und zu zählen (hier die Textsprache). Der Group-Befehl kombiniert dabei Count mit einem Zähler und Datumsgrenzen als Bedingung.

Die Abfragen mittels Count und Group stellen zwei sehr einfache und schnelle Möglichkeiten zur Häufigkeitsanalyse. Diese Methoden geben jedoch nur (gezählte) Werte in der Kommandozeile aus und erlauben weder ein direktes Speichern der Ergebnisse in Collections, noch eine Restrukturierung der Daten. Zudem sind diese Methoden nur auf 20.000 Datensätze beschränkt und bei Datensets, die mittels Sharding auf mehrere Server verteilt sind, nicht anwendbar (Trelle, 2014, S. 200). Für diese Zwecke und für komplexere Abfragen eignet sich das Aggregation Framework.

#### 4.3.2.2 Aggregation Framework

Das Aggregation Framework wurde eingeführt, um die technische Lücke zwischen den schnellen aber einfach konzipierten Abfrage-Methoden und dem komplexen aber mächtigen MapReduce zu schließen. Das Framework ist in zwei Komplexe gegliedert: einer stufenweisen *Pipeline*, die Dokumente sequentiell verarbeitet, und *Expressions*, die die Dokumente innerhalb der Pipeline manipulieren (Chodorow, 2013, S. 127-130). Eine Aggregationspipeline besteht aus einer Kette mehrerer vordefinierter Prozeduren, die sequentiell auf die Dokumente angewendet werden. So kann eine Vielzahl an Operationen mithilfe einer mehrstufigen

Pipeline verbunden werden. Die Pipeline leitet dabei die Daten von einer Prozessstufe mit jeweils spezifischen Operationen zur nächsten weiter. Abbildung 13 stellt den Ablauf der Aggregation schematisch dar.

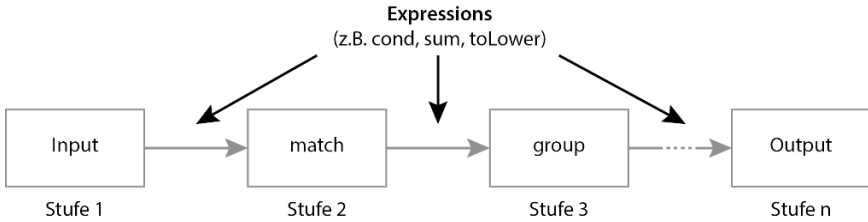


Abbildung 13: Schematische Darstellung einer Aggregation Pipeline. Eigene Darstellung.

Die Manipulation an Dokumenten erfolgt innerhalb einer Stufe der Pipeline über Expressions. Diese Expressions reagieren dabei immer nur auf den momentanen Zustand der aktuell verarbeiteten Dokumente. Es können folglich keine zusätzlichen Daten anderer Dokumente außerhalb der Pipeline hinzugefügt werden. Die möglichen Funktionen von Expressions gehen von einfachen Bedingungs-Verknüpfungen (`$and`, `$or`, ...) und Vergleichen (`$eq`, `$gt`, ...) über einfache mathematische Funktionen (`$multiply`, `$add`, ...) bis hin zur Manipulation von Strings (`$toLower`, `$toUpper`, ...).

In Listing 16 wird zunächst ein einfacher `aggregate`-Befehl ausgeführt: Alle Tweets der Collection `tatort_tweets` sollen anhand der durch Twitter erkannten Textsprache `iso_language_code` gruppiert werden. Schließlich wird die Anzahl der Tweets je Sprache gezählt und in einer absteigend sortierten Liste ausgegeben. Bereits hier zeigen sich die Vorteile gegenüber den einfachen Abfrage-Methoden: Es gibt keine Limitierung der Datenmenge<sup>27</sup> und es besteht die Möglichkeit, Ergebnisse anhand von Werten zu sortieren.

Der zweite `aggregate`-Befehl verdeutlicht das Konzept von Pipeline und Expressions. Zuerst werden Tweets ausgewählt, die im gewünschten Zeitraum veröffentlicht wurden, und anhand der Sprache gruppiert. Daneben erfasst ein Zähler die Anzahl der Tweets je Sprache und errechnet die durchschnittliche Retweet-Zahl je Sprache. Die nächste Stufe sortiert die ermittelte Liste mit den Ergebnissen absteigend anhand der Gruppengröße (nach Sprache) und begrenzt das Ergebnis

<sup>27</sup> Anmerkung: Mit zunehmender Größe der Collection steigt folglich auch die Dauer der Ausführung der Datenbank-Befehle.

auf 10 Einträge. Diese *Top 10*-Liste von Tweets und durchschnittlichen Retweets nach Sprache schreibt die letzte Stufe in die neue Collection `top`. Die nun dauerhaft gespeicherten Datensätze stehen so als Ausgangsmaterial für weitere Analysen zur Verfügung. Der `find`-Befehl gibt die Liste schließlich aus.

*Listing 16: Aggregation-Framework in MongoDB*

```
> db.tatort_tweets.aggregate(
... {$group: {_id: "$metadata.iso_language_code",
... .. count: {$sum:1}}},
... {$sort: {count: -1}});
{ "_id" : "de", "count" : 18517 }
{ "_id" : "en", "count" : 739 }
[...]
```

```
> db.tatort_tweets.aggregate(
... { $match : {"created_at":
... .. {$gte: ISODate("2015-04-11T12:00:00.000Z"),
... .. $lt: ISODate("2015-04-13T12:00:00.000Z")}}},
... { $group : {
... .. _id: "$metadata.iso_language_code",
... .. count: {$sum:1},
... .. avgRetweets: {$avg: "$retweet_count"}},
... {$sort: {count:-1}},
... {$limit: 10},
... {$out: "top"});
```

```
> db.top.find()
{ "_id" : "de", "count" : 18517, "avgRetweets" :
6.2377814980828425 }
{ "_id" : "en", "count" : 739, "avgRetweets" :
0.5940460081190798 }
{ "_id" : "und", "count" : 476, "avgRetweets" :
2.926470588235294 }
{ "_id" : "fr", "count" : 134, "avgRetweets" :
0.1417910447761194 }
[...]
```

Dieses Beispiel zeigt, dass bereits *in* der Datenbank eine sinnvolle Selektion und Umstrukturierung von Twitter-Daten erfolgen kann. So besteht die Möglichkeit, Daten anhand von (errechneten) Werten zu vergleichen, zu gruppieren und sortieren sowie die Ergebnisse in eine separate Collection zu schreiben. Jedoch gibt es

auch hier, wie bei den Abfragemethoden technische Einschränkungen: Wenn `group` und `sort`-Befehle bereits zu Beginn der Pipeline eingesetzt werden, hängt die Zahl der maximal analysierbaren Tweets von der Verfügbarkeit des Arbeitsspeichers ab. Für kumulative Operationen wie die Gruppierung oder Sortierung von Daten muss zuerst das komplette Datenset eingelesen werden. Werden dabei mehr als 10% des verfügbaren Arbeitsspeichers belegt, bricht der aggregationsprozess ab.

Neben einer guten Recherausstattung ist es daher sinnvoll, Daten in frühen Stufen der Pipeline zu filtern. Beispielsweise können einzelne Werte extrahiert und in eine neue Collection geschrieben werden. Diese Collection kann dann wiederum als Ausgangspunkt für weitere Analysen im Aggregation Framework dienen, oder – bei Export der Daten – für andere Analyse-Programme. Zudem muss beachtet werden, dass der Algorithmus zur Spracherkennung von Twitter nicht immer zuverlässig arbeitet und auch nicht alle Tweets anhand ihrer Sprache klassifizieren kann. In Listing 16 befanden sich beispielsweise 467 undefinierbare Tweets („\_id“: „,und“).

Neben dem Aggregation-Framework, das wohl für die meisten Anwendungsfälle genügt, unterstützt MongoDB noch eine weitere Technik zur Datenaggregation, die plattformübergreifend bei sehr großen und stetig wachsenden Datensets angewendet wird: *MapReduce*.

### 4.3.2.3 MapReduce

MapReduce ist eine von Google, Inc. (Dean & Ghemawat, 2008) konzipierte Technik zum kontinuierlichen Verarbeiten sehr großer Datensets. Da sie eine sehr hohe Leistungsfähigkeit hat, aber eher kompliziert und programmierintensiv in der Einrichtung ist, eignet sie sich vor allem für *Big Data*, also sehr großen und stetig wachsenden Datenmengen. Insgesamt ist die Verarbeitungsgeschwindigkeit im Vergleich zum Aggregation Framework jedoch langsamer, was auch darauf zurückzuführen ist, dass die Verarbeitung des Aggregation Frameworks vor allem im (schnellen) Arbeitsspeicher stattfindet, während MapReduce die Daten direkt auf der DB verarbeitet und die einzelnen Dokumente während der Bearbeitung (für Millisekunden) sperrt (Trelle, 2014, S. 231). Dafür können bei MapReduce deutlich mehr Daten verarbeitet werden. Der Algorithmus gliedert sich in drei Phasen: *Map*, *Group/Sort* und *Reduce*. Diese müssen zuvor einzeln mittels Javascript programmiert werden.



Jeder MapReduce-Prozess beginnt mit einer Map-Phase. Hier werden aus allen verfügbaren Dokumenten, die auf der Datenbank verteilt sind, vorher definierte Teilinformationen extrahiert. Das könnten alle Wörter eines Textes sein oder ein Zahlenwert eines bestimmten Schlüssels (wie die Anzahl der Follower eines Twitter-Nutzers). Danach folgt eine Phase des Gruppierens und Sortierens. Die extrahierten Informationen werden hier sortiert und in Bündel identischer Daten aggregiert. In der Reduce-Phase erfolgt schließlich die Reduktion der Listen-Einträge zu einem Element. (Chodorow, 2013, S. 140)

Zum besseren Verständnis zeigt das folgende, einfache Beispiel, den Ablauf des Prozesses, wie ihn Abbildung 14 vereinfacht darstellt. Ziel ist die Ermittlung der 10 häufigsten Wörter aller Tweets zum Franken-Tatort anhand der *Word Count* Methode. Das Zählen der Worthäufigkeit in Texten ist eine elementare Methode der Computerlinguistik.

Der Input für MapReduce umfasst auch in diesem Fall vollständige Tweet-Datensätze. Folglich muss der für die Analyse relevante Tweet-Text erst aus den Datensätzen extrahiert werden. Diese Erfassung aller Inhalte des Feldes `text` findet im ersten Schritt, der Map-Phase, statt. Dabei iteriert das Javascript in Listing 17 jeden Text und unterteilt diesen anhand von Wortgrenzen (hier: Leerzeichen) in einzelne Wörter. Danach folgt eine Bereinigung der Ergebnisse um Leerzeichen, Sonderzeichen und Füllwörter sowie ein Konvertieren in Kleinbuchstaben. Mit der Umwandlung in Kleinbuchstaben wird sichergestellt, dass Begriffe mit unterschiedlichen Schreibweisen (wie `TatOrt`, `tatort`, `Tatort`) als identische Werte erkannt werden können. Es entstehen nun Key-Value-Paare, in denen jedes Wort einzeln aufgelistet wird und die Häufigkeit 1 erhält, also beispielsweise `Tatort, 1`. Dabei gilt zu beachten, dass auch Wörter, die im selben Dokument mehrfach vorkommen, als einzelne Key-Value-Paare aufgelistet werden. Diese Paare tauchen in der Liste dementsprechend mehrfach als Dubletten auf und werden erst in der Reduce-Phase zusammengeführt.

In Phase 2 findet ein Sortieren und Gruppieren der Key-Value-Paare statt. Hier werden gleiche Worte dem gleichen Reducer zugewiesen um eine globale Aggregation der Wort-Häufigkeit zu ermöglichen. In diesem Beispiel erfolgen eine alphabetische Sortierung und schließlich eine Zuordnung jeweils gleicher Worte zu einem Reducer. In der anschließenden Reduce-Phase iteriert jeder einzelne Reducer über die ihm zugewiesene Wortmenge und ermittelt so die Summe der Begriffe. Aus `(Tatort, 1)` und `(Tatort, 1)` wird schließlich `(Tatort, 2)`. Am Ende des Prozesses schreiben alle Reducer ihre Ergebnisse als aggregierte Key-Value-Paare in die Ziel-Collection.

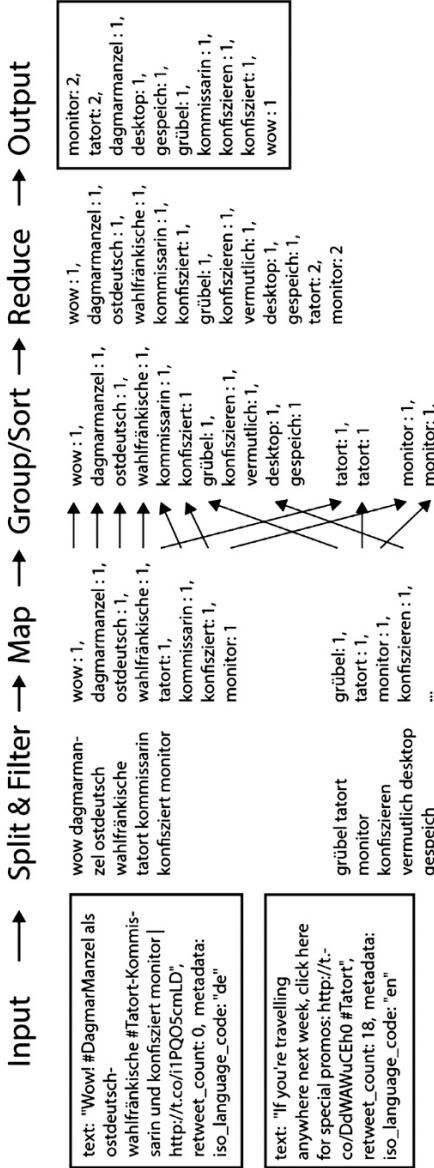


Abbildung 14: MapReduce-Prozess anhand der Word Count Methode. Eigene Darstellung.

*Listing 17:* JavaScript zur Definition der Map- und Reduce-Funktion bei Word-Count

```
// Map-Funktion
map = function() {
  if (!this.text) {
    return;
  }
  this.text.split(' ').forEach(function(word) {
    word = word.replace(/\s/g, "");
    word = word.replace(/[\.,- \/#!$%^&\*;\:;{}@=\_-'`~()]/g,
      "");
    word = word.replace(/der|die|das|ich|du|er|sie|
      es|wir|ihr|ihnen|dir|du|dem|den|
      ist|ein|eine|einer|eines|sein|mein|
      dein|im|um|auch|ja|nein|kein|
      immer|noch|und|warum|rt/gi, "");

    if(word != "") {
      emit(word.toLowerCase(), 1)
    }
  });
};

// Reduce-Funktion
reduce = function(key, values) {
  return values.length;
};
```

Der eigentliche MapReduce-Prozess wird über die MongoDB-Shell gestartet und bedient sich einer übersichtlichen Syntax. Sie beinhaltet die Namen der Map- und Reduce-Funktion sowie weitere Angaben über die Ausgabe der Ergebnisse. Wenn der jeweilige Map- und Reduce-Code bereits vorher in das Kommandofenster kopiert wurde, reicht ein einfacher Verweis auf den jeweiligen Funktionsnamen, um diesen einzubinden. Möglich wäre aber auch eine absolute oder relative Pfadangabe zum lokal gespeicherten JavaScript.

In diesem Fall speichern die Reducer ihre Ergebnisse in der Collection `wordcount`. Zuvor sortiert ein Filter alle Tweets aus, die nicht anhand der automatischen Spracherkennung von Twitter als deutschsprachig identifiziert wurden. Da der `query`-Filter vor dem Mapping einsetzt, reduziert dies unter Umständen auch die Bearbeitungszeit des Prozesses. In der Syntax können zudem weitere Verar-

beitungsschritte, wie das Errechnen von Kennzahlen (z.B. die Häufigkeit des Wortes innerhalb eines Satzes), initiiert werden. Auch besteht die Möglichkeit zum Vorsortieren und Begrenzen des Inputs für die Mapper.

*Listing 18:* Shell-Befehl zur Initiierung des MapReduce-Prozesses, nachdem die Map- und Reduce-Funktion geladen wurde

```
db.tatort_tweets.mapReduce(  
  map,  
  reduce,  
  {  
    out: "wordcount",  
    query: {"metadata.iso_language_code" : "de"}  
  }  
)
```

Jede Phase des MapReduce-Algorithmus erlaubt das Ausführen einfacher bis sehr komplexer Prozesse: Dies reicht von simplen Filtern bis zu umfangreichen Berechnungen von Kennzahlen und der Identifizierung von Mustern. Ziel des Prozesses ist dabei immer die Reduzierung des Datenumfangs auf wenige Daten oder Kennziffern. MapReduce erlaubt zudem eine hohe Flexibilität bei der Ausgabe und Speicherung der Ergebnisse: So können die erzeugten Daten in einer separaten Collection innerhalb oder außerhalb der verwendeten Datenbank abgelegt werden oder bereits vorhandene Collections überschreiben und so den Speicherbedarf verringern. Dennoch ist MapReduce für die meisten Analysen von Twitter-Daten überdimensioniert. Bereits elementare Abfragen erfordern Programmierkenntnisse in JavaScript und einen größeren Zeitaufwand als das Aggregation-Framework (Chodorow, 2013, S. 140). Folglich ist diese Analyse-methode vor allem für sehr großvolumige Berechnungen an kontinuierlich wachsenden Datensets geeignet.

#### 4.3.2.4 Vergleich der Ansätze

Mit den in MongoDB integrierten grundlegenden Abfragemethoden stehen bereits einfache Möglichkeiten zum Zählen, Filtern und Restrukturieren zur Verfügung, deren Nutzung bereits während der Datensammlung möglich ist. Für größere und schnell ansteigende Datenmengen eignen sich jedoch nur das Aggregation-Framework und MapReduce. Diese Werkzeuge ermöglichen nicht nur komplexere Ab-

fragen und Befehle, sondern funktionieren auch problemlos bei sehr großen Datensets. Allerdings ist hier die Befehls- und Prozessstruktur deutlich komplizierter und unübersichtlicher. Für MapReduce benötigt man zudem Kenntnisse in der Programmiersprache JavaScript. Im Vergleich zum Aggregation Framework ist MapReduce für sehr große Daten deutlich besser geeignet, da es keine Begrenzung hinsichtlich der zu verarbeitenden Datenmenge gibt. Allerdings ist bei dieser Methode auch der Leistungsbedarf (CPU-Auslastung) höher bei vergleichsweise geringerer Geschwindigkeit (Marturana, 2015). Tabelle 8 fasst die Vor- und Nachteile noch einmal zusammen.

Tabelle 8: Vergleich der Aggregations-Methoden von MongoDB

|                  | <b>ABFRAGEN</b>   | <b>AGGREGATION<br/>FRAMEWORK</b>  | <b>MAPREDUCE</b>   |
|------------------|---|---|--|
| <b>VORTEILE</b>  | <ul style="list-style-type: none"> <li>• Schnell</li> <li>• Einfach und übersichtlich</li> </ul>  | <ul style="list-style-type: none"> <li>• Vielzahl umfassender Funktionen</li> <li>• Stufenweiser Aufbau</li> <li>• Export der Ergebnisse</li> <li>• Sharding-Unterstützung</li> </ul> | <ul style="list-style-type: none"> <li>• Sehr mächtig</li> <li>• Kein Datenlimit</li> <li>• Export der Ergebnisse</li> <li>• Sharding-Unterstützung</li> </ul>                                   |
| <b>NACHTEILE</b> | <ul style="list-style-type: none"> <li>• Kein Speichern/Export der Ergebnisse</li> <li>• Verarbeitet maximal 20.000 Datensätze</li> <li>• Keine Sharding-Unterstützung</li> </ul> | <ul style="list-style-type: none"> <li>• Datenaufnahme bis maximal 10 % des Arbeitsspeichers</li> <li>• Programmierung unübersichtlich bei vielen Stufen</li> </ul>                   | <ul style="list-style-type: none"> <li>• JavaScript-Kenntnisse erforderlich</li> <li>• Komplex</li> <li>• Langsamer als das Aggregation Framework</li> <li>• Erhöhter Leistungsbedarf</li> </ul> |

Die in diesem Kapitel vorgestellten Methoden sind alle mit den Bordmitteln<sup>28</sup> von Python und MongoDB realisierbar. Jedoch eignen sie sich vor allem nur für allgemeine quantitative Analysen, wie der Generierung von Häufigkeitsverteilungen oder der Berechnung von Kennzahlen. Für detaillierte Inhaltsanalysen von Tweets

<sup>28</sup> Mit Ausnahme kleinerer Zusatzpakete, wie beispielsweise dem Schnittstellen-Paket *pymongo*.

bedarf es spezialisierter Software aus dem Bereich des Natural Language Processing. Das folgende Kapitel stellt deshalb ausgewählte Programme und Methoden der Computerlinguistik vor, die über das grundlegende Zählen der Worthäufigkeit aus dem letzten Beispiel hinausgehen.

### 4.3.3 *Natural Language Processing (NLP)*

Für die computergestützte Textanalyse über Python eignet sich das umfangreiche Softwarepaket *Natural Language Toolkit* (NLTK) (Loper & Bird, 2002). NLTK ist eine Sammlung von Bibliotheken und Programmen, die eine plattformunabhängige Programmierumgebung für NLP-Programme in Python bereitstellt. Es beinhaltet unter anderem standardisierte Klassen für Part-of-Speech-Tagging, Text-Klassifizierung und syntaktische Analyse (Bird et al., 2009, S. 4). Zudem besteht die Möglichkeit, Texte direkt aus einer Datenbank auf MongoDB abzurufen. NLTK wird wie andere Module installiert, benötigt aber noch zusätzliche Bibliotheken zur Darstellung von Plots oder mathematischen Berechnungen.

Die im vorigen Kapitel begonnene Analyse von Tweets zum Frankentatort wird in diesem Kapitel fortgesetzt und vertieft. Vorbereitend (siehe Listing 19) filtert ein `aggregate`-Befehl alle Tweets nichtdeutscher Sprachen außerhalb des relevanten Zeitraums und schreibt diese in eine separate Collection, die nun als Datengrundlage für den Textkorpus dient.

#### *Listing 19: Vorbereitung der NLP-Analyse in MongoDB*

```
> db.tatort_tweets.aggregate(  
... { $match : {"created_at":  
...   ... {$gte: ISODate("2015-04-11T12:00:00.000Z"),  
...   ... $lt: ISODate("2015-04-13T12:00:00.000Z")},  
...   ... "metadata.iso_language_code" : "de"}},  
... {$out: "dt_tweets"});
```

#### 4.3.3.1 Anwendungsbeispiel: Computerlinguistische Analyse des Franken-Tatorts

Zunächst wird über die MongoDB-Schnittstelle eine Verbindung zur Collection hergestellt und die Einträge ausgelesen. Als Index für diesen Textkorpus dient in diesem Fall der Veröffentlichungszeitpunkt (`created_at`). Eine Ausgabe der

Häufigkeitsverteilung aller enthaltenen Tweets über den definierten Zeitraum dient der ersten Übersicht (siehe Abbildung 15). Tweets mit dem Begriff Tatort erscheinen demnach vor allem während der Ausstrahlung und in den Stunden danach. Daraus könnte man ableiten, dass hier das typische Phänomen des *second screen* vorliegt: Zuschauer nutzen während dem Fernsehschauen weitere Geräte (wie Tablets oder Computer) und kommentieren das Gesehene in Sozialen Medien (Courtois & D'heer, 2012).

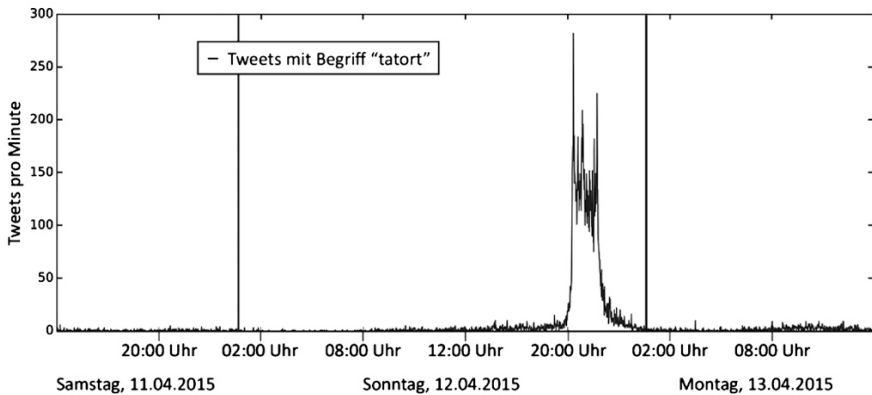


Abbildung 15: Verteilung der Tweets zum Franken-Tatort nach Uhrzeit.

Wie und über was die Twitter-Nutzer/-innen schreiben, soll nun eine detaillierte Textanalyse ermitteln. Wie bereits erwähnt wurde, bedarf die Analyse von Twitter-Text eine gründliche Vorbereitung. Im Gegensatz zu (relativ standardisierten) Buchtexten oder Nachrichtenartikeln, die eine korrekte Satzstruktur und Rechtschreibung aufweisen, bestehen Tweets häufig aus Abkürzungen, Neologismen, mehreren Sprachen, Sonderzeichen, Links und Umgangssprache (siehe dazu auch Kapitel 4.3.1). Deshalb muss der Textinput zunächst bereinigt werden, um ihn in der späteren Analyse zuverlässig zu verarbeiten.

Dafür extrahiert das Programm aus Anhang B<sup>29</sup> zunächst aus allen Tweet-Texten die einzelnen Wörter – man spricht hier von *Tokenisierung*. Um häufige Füllwörter aus der Häufigkeitsanalyse auszuschließen, sollen diese Stopwords in einem zweiten Schritt erkannt und gefiltert werden. NLTK stellt bereits Kataloge

<sup>29</sup> Aufgrund der Länge des Programmcodes erfolgt die Darstellung dieses Skript nur im Anhang.

solcher Stopwords (u.a. auf Englisch und Deutsch) zur Verfügung. Zusätzlich werden unter `customstopwords` eigene Begriffe definiert, die für den Franken-Tatort in einer hohen Häufigkeit vermutet wurden (wie zum Beispiel *Tatort*, *heute*, *gleich*). Die erkannten Wörter werden in Kleinbuchstaben umgewandelt. Das Skript liest alle Text Entities ein und filtert alle definierten Stopwords.

Zudem definiert das Skript alle Begriffe, die mit einem @-Zeichen beginnen, als Mentions und alle Wörter mit vorangestellter # als Hashtags. Alle Satzteile, die mit *http* oder *www* beginnen, werden als Links klassifiziert. Die Häufigkeitsanalyse schließt auch diese drei Wort-Klassen aus, wie auch Begriffe mit weniger als drei Buchstaben. Ebenso überprüft das Programm, ob die ermittelten Wörter rein aus Buchstaben bestehen. So fallen beispielsweise Begriffe wie *2rad* oder *kla4* aus der Untersuchung heraus. Das Histogramm in Abbildung 16 zeigt das Ergebnis der Analyse.

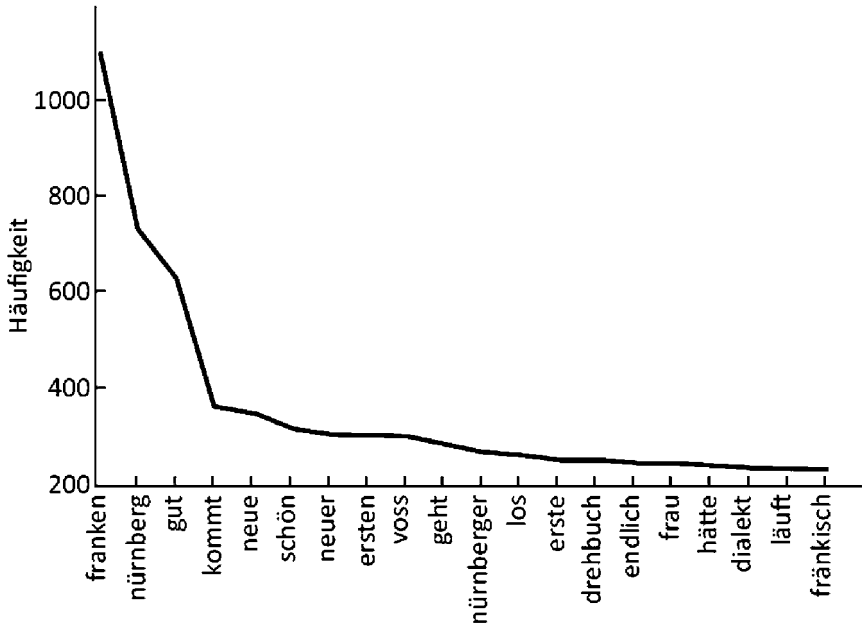


Abbildung 16: Häufigkeit der Top 20 Begriffe zum Franken-Tatort. Eigene Darstellung.



Bei Betrachtung des Resultates ragt die Prominenz des Wortes „gut“ hervor. Auf den ersten Blick scheint der Tatort eine positive Wirkung erzielt zu haben. Jedoch besteht auch die Möglichkeit, dass diese Begriffe in einem anderen Kontext verwendet wurden, wie „nicht schön“ oder „gar nicht gut“. Deswegen erscheint hier eine Betrachtung der Konkordanz, also die kontextuelle Einbettung eines Wortes in einem Satz, als sinnvoll. Für diese Auswertung wurde ein angepasstes Skript verwendet, das zwar analog zum ersten Filter alle Wörter separiert, jedoch keine Stopwords filtert (siehe Anhang B). Der Datensatz enthielte sonst beispielsweise keine wertenden Wörter wie „nicht“ oder „voll“. Das unter `sel_tokens2` gefilterte Datenset wurde nur um Links und Mentions bereinigt. Zudem erfolgt die Analyse nicht auf Wort-, sondern auf Satzebene. Auch für diese Tokenisierung stellt das NLTK eine Methode zur Verfügung. Die Separation erfolgt anhand gängiger Satzzeichen, die im Üblichen das Ende eines Satzes markieren: `.,!/?`.

#### *Listing 20:* Konkordanz des Wortes "gut" in Tweets zum Franken-Tatort

Displaying 25 of 677 matches:

```
- premiere ganz gut an . # tatort http://t.co/orng6
@ daserste sehr gut ! gute besetzung , von der handlung h
rt ich fand ihn gut mit potenzial nach oben . charaktere
rt ich fand ihn gut . solider fall und das neue duo - int
amp ; # tatort gut \ xfcbcrstanden ? - # feelixgmbh star
nderdiddle ganz gut # tatort ', u ' rt @ media_n_mngmnt :
e4mlich richtig gut ! # dadord https://t.co/z2hvg6k
der wirklich so gut ?', u '@ tatort @ brfrankentatort gib
', u ' war der gut der # tatort ? kann ihn fr \ xfcheste
rt war noch nie gut , da machen die amis bessere filme ,
ng nach richtig gut ! \ nbester spruch : do is die katz gf
rt sensationell gut ! mehr davon ', u ' ich habe kein ein
ranken - tatort gut fand :', u ' rt @ mpjhaug : hab dich
in handwerklich gut gemachter krimi viel \ u2026 ', u ' h
in handwerklich gut gemachter krimi viel besser .', u ' r
ranken - tatort gut fand :', u ' rt @ tweetbarth : viel s
er ist wirklich gut gewesen ! # megusta ', u ' top 5 der
vielleicht doch gut , dass der n \ xe4chste # bpt der # p
n n \ xfcrnberg gut ! wobei das fr \ xe4nkisch da schon e
ja ned so arch gut .', u ' rt @ loausro : @ condral902 a
# tatort ', u ' gut , dass jetzt niemand mehr \ xfcbcr de
tatort auch so gut gefallen ;) weiterhin w \ xfcschen w
haut , ist echt gut geworden :) # nuernberg # tatort ', u
ganz besonders gut hat mir am # tatort gefallen , dass n
und war richtig gut . weiter so !', u '# tatort \ n # rot
```

Der Index für das *Key Word in Context* „gut“, dargestellt in Listing 20, zeigt für die ersten Tweets eine mehrheitlich positive Einbettung des Begriffs. Um eine allgemein gültige Aussage treffen zu können, muss die Konkordanz für jeden Fall betrachtet werden. Diese Aufgabe übernimmt die Analyse von Kollokationen, die die häufigsten Wortpaare ausgibt. Das Skript im unteren Teil von Anhang B sucht nach den 20 häufigsten Bigrammen, die den Begriff „gut“ enthalten und öfter als 10 Mal im gesamten Textkorpus auftauchen. Der Output bestätigt die These der guten Bewertung des Franken-Tatorts.

#### Listing 21: Häufigste Bigramme zum Franken-Tatort

```
[(u'aufn', u'gut'), (u'gut', u'isser'), (u'richtig',
u'gut'), (u'sehr', u'gut'), (u'wirklich', u'gut'), (u'beson-
ders', u'gut'), (u'gut', u'gefallen'), (u'so', u'gut'),
(u'gut', u'tut'), (u'ganz', u'gut'), (u'ziemlich', u'gut'),
(u'gut', u'hat'), (u'gut', u'rt'), (u'echt', u'gut'),
(u'gut', u'dass'), (u'war', u'gut'), (u'gut', u'mehr'),
(u'gut', u'aber'), (u'ja', u'gut'), (u'zu', u'gut')]
```

Die bisherigen Auswertungen des Tatorts waren sehr grundlegend. Besonders bei der Ermittlung des allgemeinen Tenors auf Twitter zum Tatort genügen einfache Auszählungen oder die reine Betrachtung der Wort-Einbettung nicht. Für verlässliche Bewertungen bedarf es Algorithmen der Sentiment-Analyse, die die Stimmung beziehungsweise die Wertung von Tweets erkennen.

#### 4.3.3.2 Anwendungsbeispiel: Sentiment-Analyse von Tweets zum Franken-Tatort

Laut Pang und Lee (2008) entwickelte sich die Sentiment-Analyse zu einem mittlerweile stark repräsentierten Forschungsgebiet innerhalb der Computerlinguistik mit zahlreichen etablierten Techniken und einer Vielzahl an freien und kostenpflichtig verfügbaren Programmen. Zur automatisierten semantischen Bewertung von Tweets bedarf es immer konnotierter Korpusse, die wertende Begriffe beinhalten und diese (zum Teil) hinsichtlich ihrer Aussagekraft gewichten. Die zur Verfügung stehenden Algorithmen, wie zum Beispiel der *Naive Bayes Klassifikator*, lesen den (vorverarbeiteten) Twitter-Text ein und vergleichen die enthaltenen Begriffe mit Katalogen konnotierter Begriffe. Da die abschließende Bewertung vor allem von der Quantität und Qualität dieser Korpusse abhängt, ist die Wahl des geeigneten Korpus sehr wichtig.

Auf Twitter spezialisierte Sentiment-Korpora existieren jedoch momentan nur für die englische Sprache. Somit besteht nur die Möglichkeit, allgemeine deutsche Korpora zu nutzen oder selbst einen Korpus zu erstellen. In dieser Arbeit wird *SentiWS* (Remus, Quasthoff & Heyer, 2010) der Universität Leipzig verwendet, das kostenlos erhältlich ist. Das Programm besteht vor allem aus zwei Lexika, die 1650 positive und 1818 negativ annotierte deutsche Begriffe zuzüglich ihrer Beugungen beinhalten. Jedes Wort ist mit einem Index versehen, der von -1 (sehr negativ) bis +1 (sehr positiv) reicht.

„Ich fand den # tatort ja ned so arch gut“ (J., 2015) – bereits mit diesem beispielhaft ausgewählten Tweet wird klar, dass aufgrund der Eigenheiten der Twitter-Kommunikation nicht alle verwendeten Begriffe in diesen Lexika enthalten sein werden. Deswegen müssten falsch geschriebene Begriffe und Dialektsprache in reguläres Deutsch umgeschrieben werden. Dies könnte zum größten Teil automatisiert geschehen, sofern ein entsprechender Korpus für deutsche Dialekte und Internetjargon vorläge. Da eine manuelle Erstellung dieses Korpus den Rahmen der Arbeit überschreiten würde, greift diese nur auf den regulären SentiWS-Korpus zurück.

Basierend auf den deutschen Korpus wird im Folgenden mithilfe des überwachten maschinellen Lernens ein Identifikator trainiert, der anhand von Begriffen die emotionale Wertung eines Tweets erkennt (siehe Anhang B). Der hierfür verwendete Naive Bayes Klassifikator (NBK) nutzt die bereits vorhandenen Listen positiver und negativer konnotierter Ausdrücke des Korpus, die um zwei Listen von fiktiven Test-Tweets mit positiven und negativen Ausdrücken erweitert werden. Das Training erfolgt anhand des Korpus und wird schließlich an 3.000 zufällig ausgewählten Begriffen angewendet.

Daraus ermittelt der Naive Bayes Klassifikator 20 der aussagekräftigsten Features. Die Merkmalerkennung (*Feature Extraction*) über ein Testset dient der Verschlinkung des Ausgangswortschatzes des Klassifikators für die Gesamtanalyse. Die Reduktion auf die besonders aussagekräftigen Begriffe in Bezug auf die Stimmung eines Tweets minimiert Aufwand und Dauer der Klassifizierung. Durch die Einschränkung auf Features werden zudem auch Störfeatures, die die Qualität der Analyse minimieren, eliminiert. Das Feature Set dient nun als Wörterbuch für die Sentiment-Klassifikation des gesamten Tweet-Korpus. Nach der Erkennung der Stimmung jedes Tweets erfolgt die Ausgabe in ein Diagramm, das Abbildung 17 dargestellt.

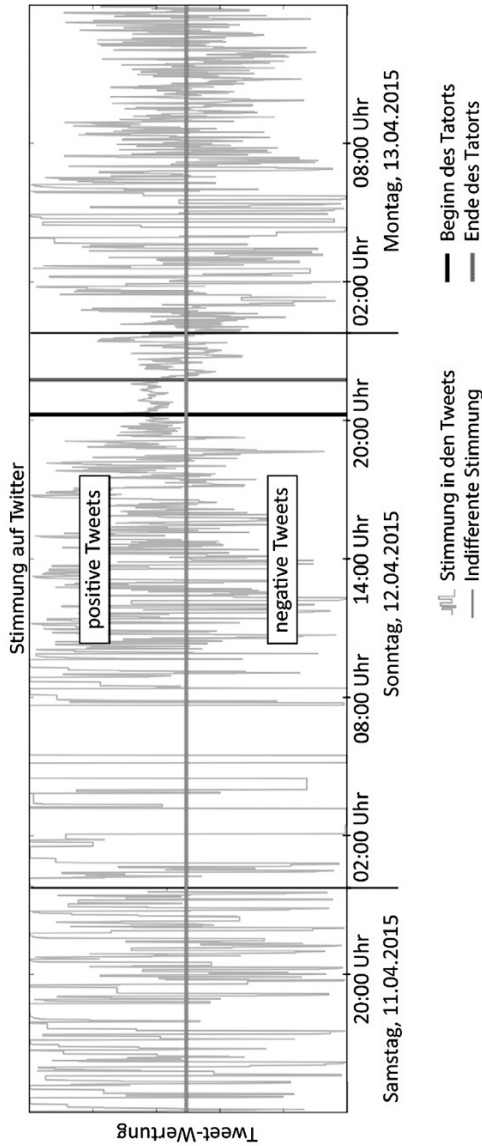


Abbildung 17: Stimmungsverlauf auf Twitter, basierend auf Tweets mit Begriff „tatort“. Eigene Darstellung.

Das Schaubild zeigt den Verlauf der aggregierten Stimmung in Tweets mit dem Begriff „tator“ über das in Listing 15 gefilterte Datenset (11. bis 13.04.). Hier wird deutlich, dass es im allgemeinen Zeitverlauf keine eindeutige Stimmung in den Tweets gibt. Dagegen erkennt man in der Zeit während der Ausstrahlung der Sendung eine überwiegend positive Stimmung. Möglich ist jedoch auch, dass der NBK schlecht trainiert wurde. Der Trainingsprozess sollte deshalb sorgfältig und idealerweise mehrstufig erfolgen: Nach Erstellung eines Trainingssets folgt die Anwendung auf ein bereits manuell konnotiertes Testset. Das Trainingsset wird dabei solange editiert, bis eine hohe Übereinstimmung zwischen manuellem und automatisch ermitteltem Sentiment vorliegt. Eine ausführliche Sentiment-Analyse über die allgemeine Darstellung der Funktionsweise hinaus übersteigt jedoch den Rahmen dieser Arbeit.

Die beiden Anwendungsbeispiele machen deutlich, dass eine verlässliche automatisierte und computergestützte Inhaltsanalyse von Tweets nicht ohne weitere Schritte der Vorverarbeitung möglich ist. Da die Qualität der Daten entscheidend für die Aussagekraft der verwendeten Algorithmen ist, sollten folglich immer vorgelagerte Filter, Strukturierungen und Korrekturen erfolgen. Dennoch besteht auch nach diesen Schritten immer noch das Problem der fehlenden Berücksichtigung des Tweet-Kontextes. Das Erkennen von Ironie und Sarkasmus ist so nicht möglich. Die Aussagekraft, besonders bei Sentiment-Analysen, ist somit immer eingeschränkt.

Nachdem nun einige Ansätze zum Verarbeiten und Analysieren von Tweets für Forschung präsentiert wurden, folgt nun eine abschließende Betrachtung von Twitter als Quelle wissenschaftlicher Arbeiten. Dabei fließen die Erkenntnisse aus den vorigen Kapiteln mit ein.

**Open Access** Dieses Kapitel wird unter der Creative Commons Namensnennung - Nicht kommerziell 4.0 International Lizenz (<http://creativecommons.org/licenses/by-nc/4.0/deed.de>) veröffentlicht, welche für nicht kommerzielle Zwecke die Nutzung, Vervielfältigung, Bearbeitung, Verbreitung und Wiedergabe in jeglichem Medium und Format erlaubt, sofern Sie den/die ursprünglichen Autor(en) und die Quelle ordnungsgemäß nennen, einen Link zur Creative Commons Lizenz beifügen und angeben, ob Änderungen vorgenommen wurden.

Erwäge Abbildungen oder sonstiges Drittmateriale unterliegen ebenfalls der genannten Creative Commons Lizenz, sofern sich aus der Abbildungslegende oder der Quellreferenz nichts anderes ergibt. Sofern solches Drittmateriale nicht unter der genannten Creative Commons Lizenz steht, ist eine Vervielfältigung, Bearbeitung oder öffentliche Wiedergabe nur mit vorheriger Zustimmung des betreffenden Rechteinhabers oder auf der Grundlage einschlägiger gesetzlicher Erlaubnisvorschriften zulässig.

