

# DSSN Tweet Corpus

Robert Rößling

University of Leipzig

2. February 2015

- 1 Task
- 2 Decision
- 3 Current state
- 4 Functioning
- 5 Demonstration
- 6 Current Problems
- 7 Solutions, Improvements and Discussion

- Creating a text corpus
  - Data should be real
  - Data should contain meta information
  - It should be possible to assign one user a certain set of data

- Creating a text corpus
  - Data should be real
  - Data should contain meta information
  - It should be possible to assign one user a certain set of data
- Should also contain social media specific attributes
  - Usernames
  - ID
  - Friends
  - Followers

- Using the Twitter REST API<sup>1</sup>
- Get real data from there
- Create User objects for each user
- Create Tweet objects for each tweet
- Give each User and Tweet a unique ID (preferably the same from Twitter)
- User object initializes one replay agent
- Possibility to create an outputfile via the native etree and BeautifulSoup4<sup>2</sup> library

---

<sup>1</sup><https://dev.twitter.com/rest/public>

<sup>2</sup><http://crummy.com/software/BeautifulSoup>

- Creates a text corpus of real data
- Date comes from Twitter
- Access through the Twitter REST API via python twitter API-tool<sup>3</sup>
- Programming Language is Python
- Every user and tweet is it's own object
- Friend requests have a unique ID
- Outputfile can be created

---

<sup>3</sup><http://mike.verdone.ca/twitter/>

<sup>4</sup><https://github.com/DSSN-Practical/DSSNTweetCorpus>

- Handler:

- Corpus:

- Friender:

- Handler:
  - Has the main method in it
  - When initiated will prompt for the initial user screen name
  - Will then prompt for the iterations
- Corpus:
  
- Friender:



- Handler:
  - Has the main method in it
  - When initiated will prompt for the initial user screen name
  - Will then prompt for the iterations
- Corpus:
  - Contains the array of all users
  - Contains an additional array of all user ID's
  - Has the authorization to the twitter API in it
- Friender:

- Handler:
  - Has the main method in it
  - When initiated will prompt for the initial user screen name
  - Will then prompt for the iterations
- Corpus:
  - Contains the array of all users
  - Contains an additional array of all user ID's
  - Has the authorization to the twitter API in it
- Friender:
  - Due to the fact that Tweets do not contain a proper "friendrequest" tweet this module is a workaround attempt
  - It generates a random subset of users for each user and sets them as "friends"

```
1 class User(object):
2     uid = 0
3     screen_name = ''
4     name = ''
5     createdAt = ''
6     description = ''
7     nrFriends = 0
8     nrFollowers = 0
9     tweets = []
10    friends = []
11    #only used for GET timeline
12    protected = False
```

```
1 class Tweet(object):
2     tid = 0
3     text = ''
4     createdAt = ''
5     isReply = False
6     replyTo = None
7     retweeted = False
8     isFriendRequest = False
9     FriendRequestToId = 0
10    #causes errors for now
11    hashtags = []
```

- startHandling

- startHandling
  - General method that will start handling all incoming data
  - Only this method needs to be called

# Handler methods

- startHandling
  - General method that will start handling all incoming data
  - Only this method needs to be called
- addUsers

- startHandling
  - General method that will start handling all incoming data
  - Only this method needs to be called
- addUsers
  - Will look up the last 200 followers of a user and will extract their meta-data
  - Creates a user object and appends it to the corpus array



- `startHandling`
  - General method that will start handling all incoming data
  - Only this method needs to be called
- `addUsers`
  - Will look up the last 200 followers of a user and will extract their meta-data
  - Creates a user object and appends it to the corpus array
- `addUserTweets`

- `startHandling`
  - General method that will start handling all incoming data
  - Only this method needs to be called
- `addUsers`
  - Will look up the last 200 followers of a user and will extract their meta-data
  - Creates a user object and appends it to the corpus array
- `addUserTweets`
  - Looks up the Timeline of a certain user
  - Creates a Tweet object for every Tweet

- `startHandling`
  - General method that will start handling all incoming data
  - Only this method needs to be called
- `addUsers`
  - Will look up the last 200 followers of a user and will extract their meta-data
  - Creates a user object and appends it to the corpus array
- `addUserTweets`
  - Looks up the Timeline of a certain user
  - Creates a Tweet object for every Tweet
- `createUserEntry` and `createOutputFile`

- `startHandling`
  - General method that will start handling all incoming data
  - Only this method needs to be called
- `addUsers`
  - Will look up the last 200 followers of a user and will extract their meta-data
  - Creates a user object and appends it to the corpus array
- `addUserTweets`
  - Looks up the Timeline of a certain user
  - Creates a Tweet object for every Tweet
- `createUserEntry` and `createOutputFile`
  - Creates HTML / XML based string
  - Saves it as a .xml File with all the user meta-data and tweets

# Demonstration

- Friend request data is not "real"

- Friend request data is not "real"
- User can send himself a friend request

- Friend request data is not "real"
- User can send himself a friend request
- Creating the outputfile results in a high RAM usage



- Friend request data is not "real"
- User can send himself a friend request
- Creating the outputfile results in a high RAM usage
- Creating the outputfile on Windows results in a crash

- Friend request data is not "real"
- User can send himself a friend request
- Creating the outputfile results in a high RAM usage
- Creating the outputfile on Windows results in a crash
- Halt method slows down the computation drastically

- Use the lxml<sup>5</sup> library for high-performance XML parsing

---

<sup>5</sup><http://lxml.de/>

# Solutions, Improvements and Discussion

- Use the lxml<sup>5</sup> library for high-performance XML parsing
- Return an user object to the replay agent directly for easy handling

---

<sup>5</sup><http://lxml.de/>

# Solutions, Improvements and Discussion

- Use the lxml<sup>5</sup> library for high-performance XML parsing
- Return an user object to the replay agent directly for easy handling
- From the actual amount of friends of each user generate a more "real" subset

---

<sup>5</sup><http://lxml.de/>

- Use the lxml<sup>5</sup> library for high-performance XML parsing
- Return an user object to the replay agent directly for easy handling
- From the actual amount of friends of each user generate a more "real" subset
- Parallel halt method:
  - Assumption: Friends = Followers, since only Friend requests are used
  - While Rate limit for GET/followers/list is exceeded use GET/friends/list
  - while both Rate limits would be exceeded already start reading out the user timelines
  - **Problem:** Amount of users increases faster while followers / friends request are linear, which results into several thousand user timeline requests for only 20+ users, the waiting time for user timeline can hardly be reduced

---

<sup>5</sup><http://lxml.de/>