



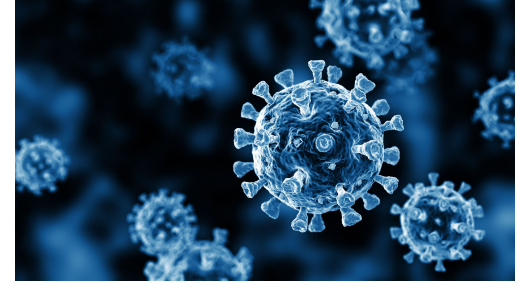
Geological Extraction for Twitter



- Varad Pimpalkhute, affiliated with Indian Institute of Information Technology, Nagpur.
- Project Advisor: Dr. Arjun Magge.
- Supervisor: Dr. Graciela Gonzalez-Hernandez.

Why do we need geolocation of a social media user?

- Tracking Infectious Diseases.
- Gathering analysis on how well an advertisement is received in a particular region.
- Keeping tabs on spread of diseases such as COVID-19.



Task : Given an input tweet and user screen name, identify the user's current location of residence.

Let's Take an Example

User Screen Name :



A screenshot of a Twitter profile for Chelsea Finn. The profile picture shows a woman with dark hair smiling. The name is **Chelsea Finn** with the handle [@chelseabfinn](#). The bio reads: "CS Faculty @Stanford. Research scientist @GoogleAI. PhD from @Berkeley_EECS, EECS BS from @MIT". There is a [#BlackLivesMatter](#) hashtag and a location of Palo Alto, CA. The profile shows 310 following and 35.8K followers. A small video thumbnail is visible in the bio area.

Input User Tweet :



A screenshot of a tweet from Chelsea Finn (@chelseabfinn). The tweet text is: "Last week, I gave a talk 'at' Toronto discussing 3 principles for tackling distribution shift: * pessimism * adaptation * anticipation and how they can improve robustness to spurious correlations, changes in users, and non-stationary & multi-agent RL envs." Below the text is a video player showing a slide titled "Principles for Tackling Distribution Shift: Pessimism, Adaptati...". The video player includes a play button and a URL: [youtube.com](#). The tweet is dated 5:56 AM · Mar 2, 2021 · Twitter Web App. The engagement statistics are 48 Retweets, 1 Quote Tweet, and 286 Likes. The bottom of the tweet shows icons for replying, retweeting, liking, and sharing.

A few assumptions and challenges

- Old datasets are publicly available for use.
- We assume that a given user has not changed his place of residence.
- If an account is deleted, extracting metadata becomes arduous.
- We assume that user feeded information is true.
- Tweet with metadata amounts to only 1% of all the tweets available.

```
{
  "created_at": "Sat Jan 10 17:18:56 +0000 2009",
  "default_profile": false,
  "default_profile_image": false,
  "description": "Prime Minister of India",
  "favourites_count": 0,
  "follow_request_sent": false,
  "followers_count": 66660683,
  "following": true,
  "friends_count": 2354,
  "geo_enabled": false,
  "has_extended_profile": true,
  "id": 18839785,
  "id_str": "18839785",
  "is_translation_enabled": false,
  "is_translator": false,
  "lang": null,
  "listed_count": 26702,
  "location": "India",
  "name": "Narendra Modi",
  "notifications": false,
  "profile_location": null,
  "profile_sidebar_border_color": "FFFFFF",
  "profile_sidebar_fill_color": "D5DFED",
  "profile_text_color": "233863",
  "profile_use_background_image": true,
  "protected": false,
  "screen_name": "narendramodi",
  "status": {
    "contributors": null,
    "coordinates": null,
    "created_at": "Tue Mar 30 17:17:06 +0000 2021",
    "entities": {
      "hashtags": [],
      "media": [

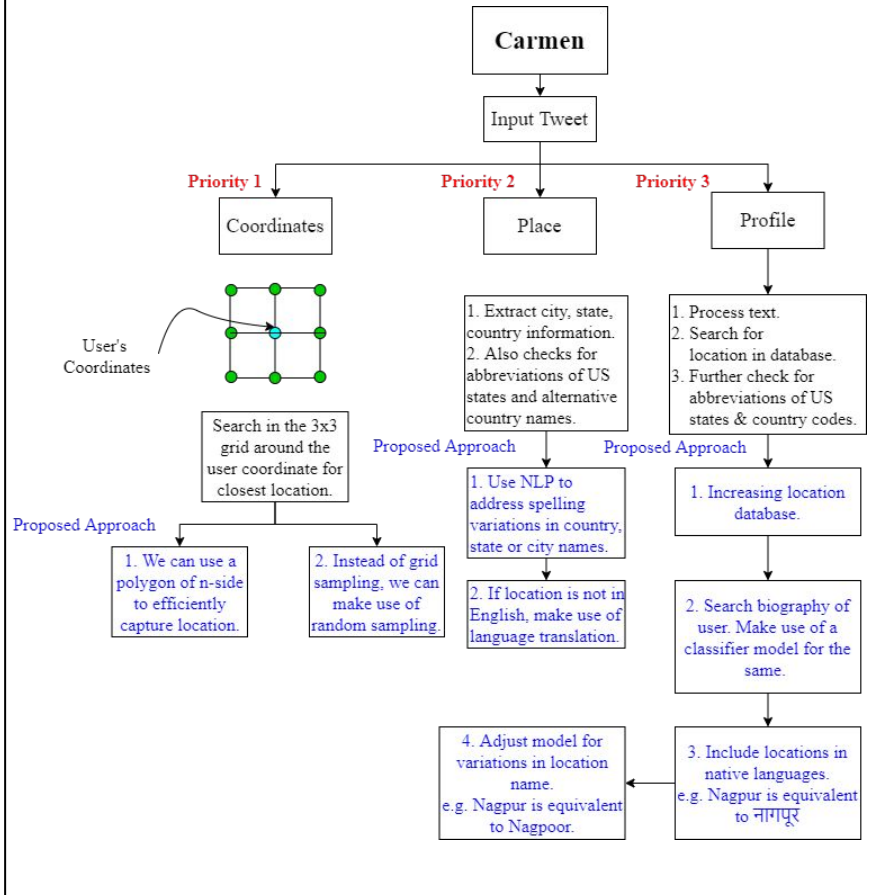
```

USER METADATA

Some of the important location related tags are:

- ★ Geo tag : (latitude, longitude)
- ★ Place tag : Location based on coordinates
- ★ Profile tag : Location form present in profile.

Carmen

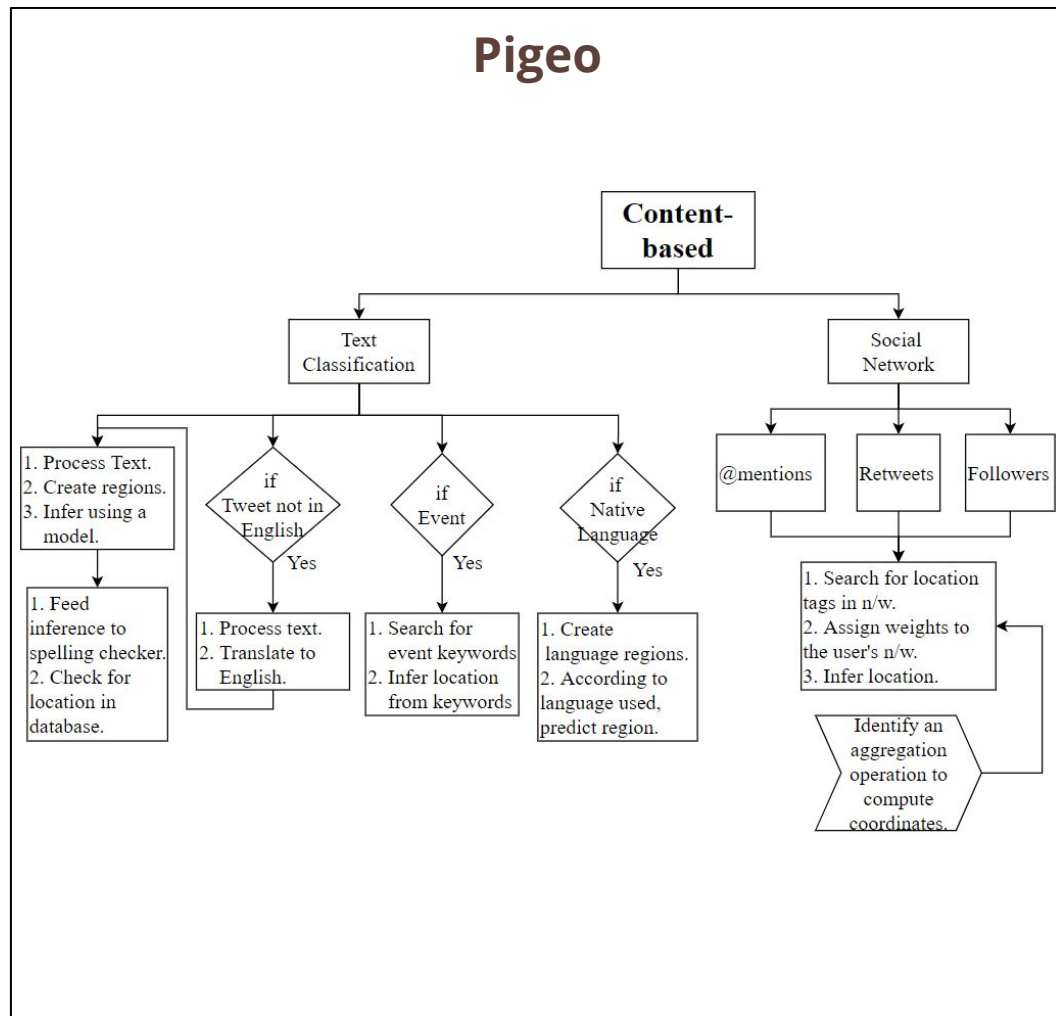


Exemplar methods in literature

- + Simple and fast
- + Tends to cover all meta tags
- Small location database
- Doesn't account for content-based information.
- Prone to errors due to spelling variations.

Exemplar methods in literature

- + Accounts for content-based information.
- + Makes use of social network of a Twitter user.
- Doesn't account for meta information.
- Accuracy is low.
- Doesn't structure the location hierarchically.



Goal

1. Improve on Carmen and Pigeo. Finally, integrate both the methodologies in one model.
2. Generating a dataset for country-wise classification by scrapping user data from Twitter using tools like Tweepy.
3. Addressing the issue of classification on a hierarchical level.



Methodology

Content-based Approach

- Generating country-wise dataset for classification.
- Building different models for text classification using FastText, Transformers, CNNs, etc.
- Fine tuning models.
- Addressing the issue of hierarchical classification.

Metadata-based Approach

- Generating a location database.
- Introducing population heuristics for better prediction of location.
- Computing the radius of each location using GADM and geopandas.

➤ Generating country-wise dataset for classification.

```
1 # UK: 12 regions
2 uk_states = ['Scotland', 'Wales', 'North East', ..., 'West Midlands', 'East Midlands', 'Greater London']
3
4 # main script
5 user_list, tweets = [], []
6 for state in states:
7     user_state_list = []
8     # Searching for users living in the state by making use of 4 different keywords and updating the
9     # state dictionary accordingly
10    user_state_list = keyword_search(state, 1, user_state_list)
11    user_state_list = keyword_search(state, 2, user_state_list)
12    user_state_list = keyword_search(state, 3, user_state_list)
13    user_state_list = keyword_search(state, 4, user_state_list)
14    # Update the main dictionary
15    user_list.extend(user_state_list) # or we can do something like this:
16    # user_dict[state] = user_state_dict
17    # Extracting Tweets from the user screen name.
18    for user in user_state_list:
19        tweet = user_tweets(state, user)
20        if tweet is not None:
21            tweets.extend(tweet)
22    print('Completed Scrapping', state)
```

➤ Generating country-wise dataset for classification.

```
Completed Scrapping Scotland
Completed Scrapping Northern Ireland
Completed Scrapping Wales
Completed Scrapping North East
Completed Scrapping North West
Completed Scrapping Yorkshire and the Humber
Completed Scrapping West Midlands
Completed Scrapping East Midlands
Completed Scrapping South West
User has protected tweets, skipping ...
Completed Scrapping South East
Completed Scrapping East of England
Completed Scrapping Greater London
Total Tweets Scrapped: 122152
Total user IDs: 1272

(geoloc) C:\Users\Varad\OneDrive\Desktop>python create_dataset.py
```

- Building different models for text classification using
 - FastText,
 - Transformers,
 - CNNs, etc.
- Generating a location database.
- Computing the radius of each location using GADM and geopandas.

Future work

- ★ Improve on the classification.
- ★ Develop a hierarchical classifier.
- ★ Make use of social network of a Twitter user.
- ★ Account for information in native languages.

References

- [1] X. Zheng, J. Han and A. Sun, "A Survey of Location Prediction on Twitter," in IEEE Transactions on Knowledge and Data Engineering, vol. 30, no. 9, pp. 1652-1671, 1 Sept. 2018, doi: 10.1109/TKDE.2018.2807840.
- [2] Han, T. (2012). Geolocation Prediction in Social Media Data by Finding Location Indicative Words. In Proceedings of COLING 2012 (pp. 1045–1062). The COLING 2012 Organizing Committee.
- [3] Roller, J. (2012). Supervised Text-based Geolocation Using Language Models on an Adaptive Grid. In Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (pp. 1500–1510). Association for Computational Linguistics.
- [4] Roesslein, J. (2020). Tweepy: Twitter for Python!. URL: <https://github.com/tweepy/tweepy>.
- [5] Mark Dredze, Michael J Paul, Shane Bergsma, & Hieu Tran (2013). Carmen: A Twitter Geolocation System with Applications to Public Health. In AAAI Workshop on Expanding the Boundaries of Health Informatics Using AI (HIAI).
- [6] Rahimi, A., Cohn, T., & Baldwin, T. (2016). pigeo: A Python Geotagging Tool. In Proceedings of ACL-2016 System Demonstrations (pp. 127–132). Association for Computational Linguistics.