

Coping with GDPR in Social Media Research



Universitat
Pompeu Fabra
Barcelona

Marc Vives
Data Protection Officer
Universitat Pompeu Fabra
dpo@upf.edu



GDPR: General Data Protection Regulation.
Regulation (EU) 2016/679 on the protection of natural persons with regard to the processing of their personal data

Personal data: Any information related to an identified or identifiable individual (the *data subject*)



- Personal data is pervasive in social media
- We will use Twitter as an example
- All Social Media platforms share similar situations



What's in a Tweet?

Twitter API v2
data dictionary

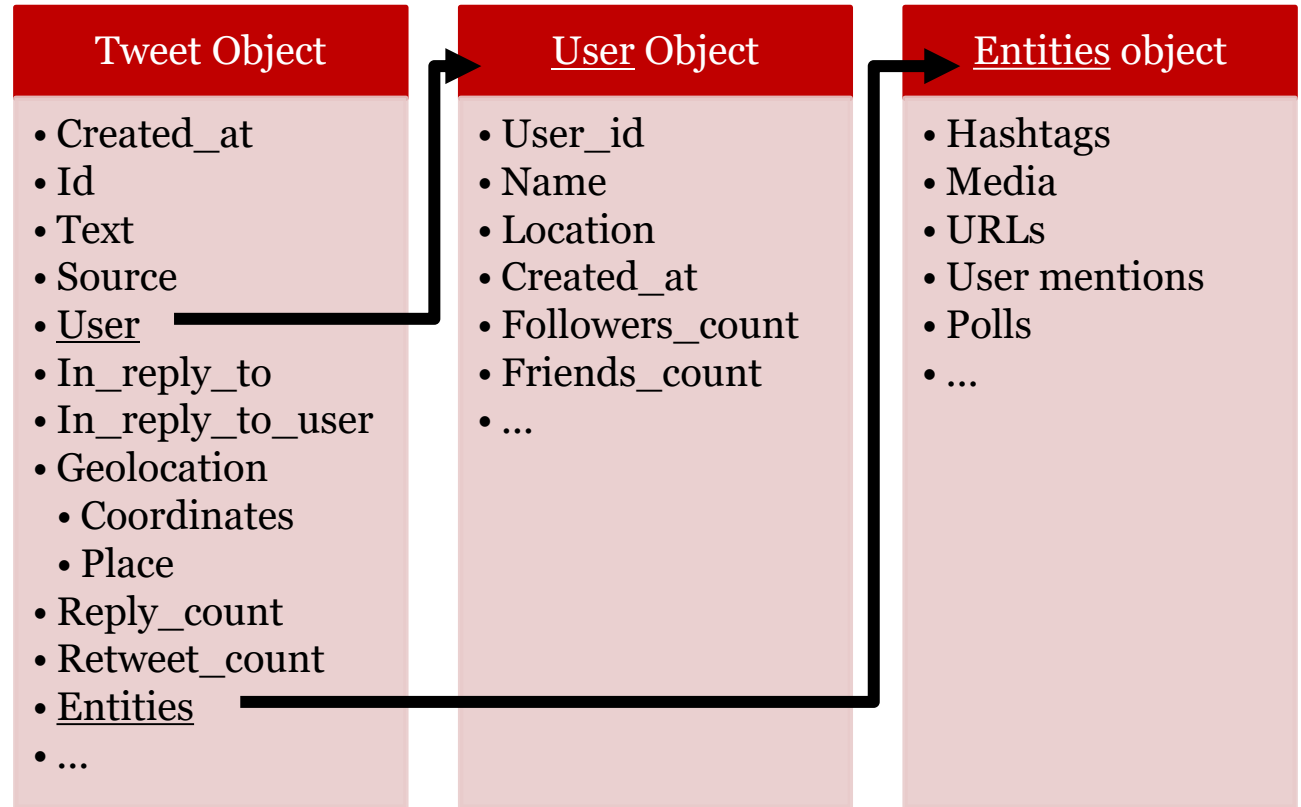
Tweet Object

- Created_at
- Id
- Text
- Source
- User
- In_reply_to
- In_reply_to_user
- Geolocation
 - Coordinates
 - Place
- Reply_count
- Retweet_count
- Entities
- ...



What's in a Tweet?

Twitter API v2
data dictionary





What's in a Tweet?

Data related to identified data subjects

Tweet Object
<ul style="list-style-type: none">• Created_at• Id• Text• Source• User• In_reply_to• In_reply_to_user• Geolocation<ul style="list-style-type: none">• Coordinates• Place• Reply_count• Retweet_count• Entities• ...

User Object
<ul style="list-style-type: none">• User_id• User_Name• Location• Created_at• Followers_count• Friends_count• ...

Entities object
<ul style="list-style-type: none">• Hashtags• Media• URLs• User mentions• Polls• ...



What's in a Tweet?

Data related to identifiable data subjects



Tweet Object
• Created_at
• Id
• Text
• Source
• <u>User</u>
• In_reply_to
• In_reply_to_user
• Geolocation
• Coordinates
• Place
• Reply_count
• Retweet_count
• <u>Entities</u>
• ...

<u>User</u> Object
• User_id
• Name
• Location
• Created_at
• Followers_count
• Friends_count
• ...

<u>Entities</u> object
• Hashtags
• Media
• URLs
• User mentions
• Polls
• ...



What's in a Tweet?

Data profiling data subjects (affinities, habits, ideology...)

Tweet Object	User Object	Entities object
<ul style="list-style-type: none">• Created_at• Id• Text• Source• User• In_reply_to• In_reply_to_user• Geolocation<ul style="list-style-type: none">• Coordinates• Place• Reply_count• Retweet_count• <u>Entities</u>• ...	<ul style="list-style-type: none">• User_id• Name• Location• Created_at• Followers_count• Friends_count• ...	<ul style="list-style-type: none">• Hashtags• Media• URLs• User mentions• Polls• ...



GDPR stipulates that any processing of personal data shall be grounded in one of 6 possible legal bases.

In the vast majority of cases the legal basis for processing personal data in a research project is data subjects' consent

So, what can we do to conduct research with personal data extracted from social media networks?

- 1) If possible, seek users' consent to process their data
- 2) If not possible, anonymise the data prior to their storage and processing

Anonymisation: application of a process to personal data so the link to the data subject is permanently removed

GDPR is not applicable to anonymised data



- The rest of this presentation is about techniques to anonymise data collected from social media



Coping with text

Tweet Object
• Created_at
• Id
• Text
• Source
• User
• In_reply_to
• In_reply_to_user
• Geolocation <ul style="list-style-type: none">• Coordinates• Place
• Reply_count
• Retweet_count
• Entities
• ...

User Object
• User_id
• User_Name
• Location
• Created_at
• Followers_count
• Friends_count
• ...

Entities object
• Hashtags
• Media
• URLs
• User mentions
• Polls
• ...



Coping with text

1. Keyword extraction

INPUT

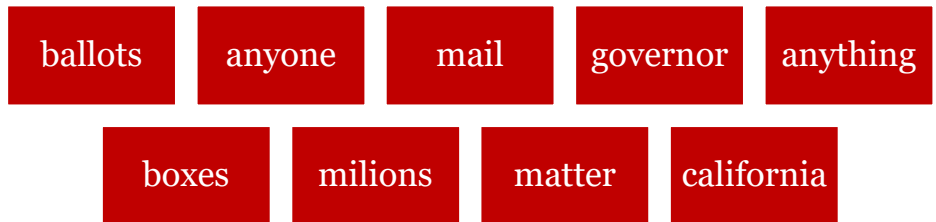


Donald J. Trump
@realDonaldTrump

There is NO WAY (ZERO!) that Mail-In Ballots will be anything less than substantially fraudulent. Mail boxes will be robbed, ballots will be forged & even illegally printed out & fraudulently signed. The Governor of California is sending Ballots to millions of people, anyone.....

[Traducir Tweet](#)

OUTPUT





Coping with text

2. Translation

INPUT



Donald J. Trump
@realDonaldTrump



There is NO WAY (ZERO!) that Mail-In Ballots will be anything less than substantially fraudulent. Mail boxes will be robbed, ballots will be forged & even illegally printed out & fraudulently signed. The Governor of California is sending Ballots to millions of people, anyone.....

[Traducir Tweet](#)

OUTPUT

No hay ninguna posibilidad (icero!) de que los votos por correo sean algo menos que sustancialmente fraudulentos. Los buzones serán robados, las papeletas serán falsificadas e incluso impresas ilegalmente y firmadas fraudulentamente. El Gobernador de California está enviando papeletas a millones de personas, cualquier persona...

<https://www.DeepL.com/Translator>



Coping with text

3. Other tricks

Remove empty words (with no lexical meaning)	<i>Mail boxes will be robbed, ballots will be forged & even illegally printed out & fraudulently signed.</i>
Remove direct identifiers	<i>The Governor of California is sending ballots to millions of people</i>
Replace derivative words with their root form	<i>ballots → ballot will be → be fraudulent → fraud</i>
Replace words with a synonym	<i>NO alternative ZERO Mail-In Ballots is anything less substantially dishonest.</i>
Use a combination of techniques	Best results always obtained using a combination of techniques!!!



Coping with images

Tweet Object	User Object	Entities object
<ul style="list-style-type: none">• Created_at• Id• Text• Source• User• In_reply_to• In_reply_to_user• Geolocation<ul style="list-style-type: none">• Coordinates• Place• Reply_count• Retweet_count• Entities• ...	<ul style="list-style-type: none">• User_id• User_Name• Location• Created_at• Followers_count• Friends_count• ...	<ul style="list-style-type: none">• Hashtags• Media• URLs• User mentions• Polls• ...



Coping with images



<https://docs.imagg.com/#getting-started-request>

Semantic extraction

```

{ "result":
  {
    "tags": [
      {
        "confidence": 61.4116096496582,
        "tag": {
          "en": "mountain"
        }
      },
      {
        "confidence": 54.3507270812988,
        "tag": {
          "en": "landscape"
        }
      },
      {
        "confidence": 50.969783782959,
        "tag": {
          "en": "mountains"
        }
      },
      {
        "confidence": 46.1385192871094,
        "tag": {
          "en": "wall"
        }
      },
      {
        "confidence": 16.6996059417725,
        "tag": {
          "en": "lake"
        }
      },
      {
        "confidence": 16.6136302947998,
        "tag": {
          "en": "cliff"
        }
      },
      {
        "confidence": 16.5426540374756,
        "tag": {
          "en": "geology"
        }
      },
      {
        "confidence": 15.9809865951538,
        "tag": {
          "en": "wilderness"
        }
      },
      {
        "confidence": 9.03097343444824,
        "tag": {
          "en": "roof"
        }
      },
      {
        "confidence": 8.87552165985107,
        "tag": {
          "en": "peaks"
        }
      },
      {
        "confidence": 8.81966876983643,
        "tag": {
          "en": "alpine"
        }
      },
      {
        "confidence": 8.80224514007568,
        "tag": {
          "en": "mount"
        }
      },
      {
        "confidence": 8.73800754547119,
        "tag": {
          "en": "f"
        }
      }
    ]
  }
}

```




Coping with dates and numbers

Tweet Object
• Created_at
• Id
• Text
• Source
• User
• In_reply_to
• In_reply_to_user
• Geolocation
• Coordinates
• Place
• Reply_count
• Retweet_count
• Entities
• ...

User Object
• User_id
• User_Name
• Location
• Created_at
• Followers_count
• Friends_count
• ...

Entities object
• Hashtags
• Media
• URLs
• User mentions
• Polls
• ...



Coping with dates and numbers

Anonymization by aggregation	<i>Use statistical values (mode, median, standard deviation..) instead of keeping individual figures.</i>
Adding white noise (random number)	<i>Add $X, X \sim U(-a, a)$ to each individual figure, so statistical properties of the sample will not change.</i>
Use of hash algorithms	<i>Hash algorithms are susceptible to brute force attacks. 'Salted' hash minimizes the issue.</i>
Remove a random number of samples	<i>If there is certainty that an element is in the sample, it is easier to find.</i>
Use data in ranges	<i>E.g., if you have to process account creation date, use ["< one year old", "between 1 and 5", "more than 5 y.o."]</i>
Always perform a reidentification test	<i>Best results if reidentification test is conducted by a colleague</i>



To sum up:

- If you make use of personal data extracted from social media for your research:
 1. Seek users' consent, if possible
 2. Anonymise the data by adapting the ideas expressed in this presentation to your research and
 - Do not store interim personal data
 - Anonymisation shall be conducted 'on the fly'
 - If you need to recreate the dataset, re-run the process
 3. Always seek certification by a research ethics committee
 - Research with social media data is an evolving topic
 - Ideas expressed in this presentation may change in the future



**Universitat
Pompeu Fabra**
Barcelona

GDPR stipulates that any processing of personal data shall be grounded in one of 6 possible legal bases.

In the vast majority of cases the legal basis for processing personal data in a research project is data subjects' consent. But GDPR provides extra room for the so called 'further processing': when personal data is used for a processing different than the one for which data was collected

Considerations when further processing of personal data (GDPR's art. 6.4)

Social Media users posted their personal data for a purpose, researchers 'further process' that data for a different purpose (research)

- 1) Linking between initial purpose and intended further processing
- 2) Context in which the personal data was collected
- 3) Nature of personal data (special categories of PD, profiling data...)
- 4) Possible consequences of the further processing
- 5) Existence of appropriate safeguards