

Project Report

Introduction

This project requires us to crawl the posts, photos and corresponding comments of two users, **Swisher Sweets** and **Backwoods**, on Instagram from 2018 to 2020. After that, annotate the each posts and comments, which are text format, with 3 different labels, which is whether they contain tobacco use warnings or not, whether they contain explicit sponsorship hashtags, and whether they contain ambiguous sponsorship hashtags.

Crawling Part

Details

1. Before crawling, we first need to parse the URL of the corresponding user's instagram homepage. And URL of the user's instagram homepage is the URL given in the project requirement, like '<https://www.instagram.com/swishersweets/>'. We can obtain the text of the html file by sending a get request to this URL. After getting the source code of the html, we find that all the posts pages or URLs is in a javascript object.

```
283 <link rel="stylesheet" href="/static/bundles/es6/ConsumerUICommons.css/8497b02914cc.css" type="text/css" crossorigin="anonymous" />
284 <link rel="stylesheet" href="/static/bundles/es6/Consumer.css/57e5336fb32d.css" type="text/css" crossorigin="anonymous" />
285 <script type="text/javascript">window._sharedData = {"config":{"csrf_token":"bkmllyx25ncc3hp83lykpjyCvXjXb","viewer":{"biography":"","category_name":null,"exter
286 <script type="text/javascript">window._initialDataLoaded(window._sharedData);</script>
287 <script type="text/javascript">var __BUNDLE_START_TIME__=this.nativePerformanceNow?nativePerformanceNow():Date.now(),__DEV__=false,process=this.process||{};process
288 __s={"js":{"149":"/static/bundles/es6/PasswordEncryptionLogger.js/a6100073b01d.js","150":"/static/bundles/es6/EncryptionUtils.js/c5ce5fb96f3c.js","151":"/static/bu
289 <script type="text/javascript" src="/static/bundles/es6/Vendor.js/48e0f28aa478.js" crossorigin="anonymous"></script>
290 <script type="text/javascript" src="/static/bundles/es6/zh_cn.js/8da79405c5c5.js" crossorigin="anonymous"></script>
291 <script type="text/javascript" src="/static/bundles/es6/ConsumerLibCommons.js/e21525839a5f.js" crossorigin="anonymous"></script>
292 <script type="text/javascript" src="/static/bundles/es6/ConsumerUICommons.js/ecc59c78c1fd.js" crossorigin="anonymous"></script>
293 <script type="text/javascript" src="/static/bundles/es6/ConsumerAsyncCommons.js/c4ca4238a0b9.js" crossorigin="anonymous"></script>
294 <script type="text/javascript" src="/static/bundles/es6/Consumer.js/d575aab3c2da.js" crossorigin="anonymous" charset="utf-8" async=""></script>
295 <script type="text/javascript" src="/static/bundles/es6/ProfilePageContainer.js/7c0a1a1d1979.js" crossorigin="anonymous" charset="utf-8" async=""></script>
```

We can get the first 12 posts urls in this object, but the first 12 posts is created in 2021, and what we need is the posts from 2018 to 2020.

2. At this point, we need to find how to find more posts in the homepage. When we scroll down the page, more pages occur, and new requests in the browser are generated. We find that the browser send request to this page.

```
Request URL: https://www.instagram.com/graphql/query/?query_hash=32b14723a678bd4628d70c1f877b94c9&variables=%7B%22id%22%3A%221717894442%22%2C%22first%22%3A12%2C%22after%22%3A%22QVFBY0Fld2NiX1RkOFRCpCHVQOUo2SGpLVlFLZU1KV0hMeUk1dVpZY2h4OWV5Q2RmYjZtbXhtcm85TTNIa2xwTVdEdWctN1VFa3d1c19wMlVBOWdFREFKVA%3D%3D%22%7D
```

And we look up the response of this request. We can see that there are `has_next_page` and `end_cursor` two values, and actually we can get next 12 posts by using `end_cursor` and the `user_id`, which got from the source page.

```
{
  "data": {
    "user": {
      "edge_owner_to_timeline_media": {
        "count": 2114,
        "page_info": {
          "has_next_page": true,
          "end_cursor": "QVFBY0Fld2NiX1RkOFRCpCHVQOUo2SGpLVlFLZU1KV0hMeUk1dVpZY2h4OWV5Q2RmYjZtbXhtcm85TTNIa2xwTVdEdWctN1VFa3d1c19wMlVBOWdFREFKVA=="
        }
      }
    }
  },
  "status": "OK"
}
```

3. And we find that all the data, like photos URLs, posts and time to create the posts are contained in a json object.

```

▼ edges: [(node: {__typename: "GraphImage", id: "2528716879783823442", gating_info: null,...}),...]
▼ 0: {node: {__typename: "GraphImage", id: "2528716879783823442", gating_info: null,...}}
  ▼ node: {__typename: "GraphImage", id: "2528716879783823442", gating_info: null,...}
    accessibility_caption: null
    comments_disabled: false
    dimensions: {height: 1080, width: 1080}
    display_resources: [(...), (...), (...)]
    display_url: https://scontent-tpe1-1.cdninstagram.com/v/t51.2885-15/e35/159444221_116996823777225_2861468624840595764_n.jpg?tp=1&nc_ht=scontent-tpe1-1.cdninstagram.com&nc
    edge_media_preview_likes: {count: 662, edges: []}
    edge_media_to_caption: {edges: [(node: {text: "It's about that time..."})]}
    edge_media_to_comment: {count: 29, page_info: {has_next_page: true, end_cursor: ""}}
    edge_media_to_sponsor_user: {edges: []}
    edge_media_to_tagged_user: {...}
    fact_check_information: null
    fact_check_overall_rating: null
    gating_info: null
    id: "2528716879783823442"
    is_paid_partnership: false
    is_video: false
    location: null
    media_overlay_info: null
    media_preview: "ACoqyUPWryD5as7Q3UZqNolZ7enattkVYvFAZ6/ugZOP6fU1XL7+VAHQ0oP8+vuIptFVLBepuFpm54DC5HXoazYv8yEEjkhBqXzV7rzTGPFAQ0kwlM0PyDvVaSuMttu/LjDYx+hp0LYUWVljCjK9CT+MeP6i
    owner: {id: "1717894442", username: "backwoods_cigars"}
    sensitivity_friction_info: null
    sharing_friction_info: {should_have_sharing_friction: false, bloks_app_url: null}
    shortcode: "CkX2apq90Hs"
    taken_at_timestamp: 161566564

```

We just need to save what we need in a list to build our own dataset, and when searching for each post, we should compare the `taken_at_timestamp` with the start timestamp and the end timestamp to filter out the posts from 2018 to 2020.

4. According to the comments, all the comments are saved in an independent URL, same as posts.
5. After saving what we need, such as photos URLs, comments and posts, we just write the comments and posts on disks, then iterate all the photos URLs to send a get request and download the photos.

Problems

It is tough to crawl these data, since we have not applied the open-source tools, instead, we analyze the source code of the html, and the logic behind the request.

1. When we crawl, we actual do a request to the URL. And to request to the instagram, there should be a certification process. At this point, we should copy the cookie in the browser to the request header.
2. It is easy to reach the instagram's connection request limit, we don't record what the limit is, but if we send request to instagram many times in a short time, it will return a 429 code. So, we add a sleep action after each request action, when we receive a 429 code, we sleep a time and retry to send request.

Annotation Part

Explicit and Ambiguous Sponsorship Hashtags

Simple. We decide whether the sentence, among posts and comments, contains (#sponsored, #ad, #paid, #thanks, #sp, #spon, #ambassador or #collab).

Tobacco Use Warnings

In this part, we do not train a model, because we do not have enough labeled samples to build a training set. So we find some example sentences in the following pages.

- <https://www.fda.gov/tobacco-products/labeling-and-warning-statements-tobacco-products/cigarette-labeling-and-health-warning-requirements>
- https://www.cdc.gov/tobacco/data_statistics/sgr/2000/highlights/labels/index.htm

After that, we use sentence embedding to decide whether the posts and the comments are similar to example sentences, and the embedding vectors are downloaded from <https://tfhub.dev/google/LaBSE/1>.

We use the sample code from the page to do embedding. While experimenting, it's easy to reach the limit of RAM, since when we send many sentences together to the model, the weights size can be extremely large. So like the comments, we set a batch size as 1000, and the max sentence length 256 to do embedding.

After embedding finishing, we just multiply the comments embeddings or the posts embeddings with the example embeddings, note that this is matrix multiplication. And we find the new matrix line by line, if there is a value higher than the similarity threshold we set manually, we annotate that sentence contain tobacco use warnings.