

Crawlers:

Each forum webpage contains multiple post URLs. The multiprocessing crawlers extract post URLs from the forum webpages and crawl them. The crawled contents are stored by using MongoDB. Failed URLs will be retried up to three times.

Data format:

```
_id: ObjectId("612a1211cbb1f357df0d1083")
title: "Dns error - Spiceworks"
url: "https://community.spiceworks.com/topic/2785-dns-error?from_forum=215"
closed: true
user_list: Array
  0: "Elias"
  1: "Marcelo (Spiceworks)"
original_post: ""IP Address does not resolve to a hostname"
               ping and nslookup up both r..."
created_date: "May 8, 2007 at 14:59 UTC"
post_creator: "Elias"
replies: Object
  Marcelo (Spiceworks): Array
    0: Object
        date: "May 8, 2007 at 16:21 UTC"
        reply: "A little more clarification please. Are you getting the error before o..."
    1: Object
        date: "May 9, 2007 at 11:28 UTC"
        reply: "With the new computer name try removing the workstation from the domai..."
    Elias: Array
      0: Object
          date: "May 9, 2007 at 11:18 UTC"
          reply: "I changed the Computer name on the computer and then I got the error. ..."
      1: Object
          date: "May 9, 2007 at 11:45 UTC"
          reply: "ok, I'll try it . However, IF I want to change my workstation naming s..."
  external_urls: Array
```

user_list: stores users who have replied to the post.

original_post: the original question / feedback wrote by the post creator.

replies: replies made by each user are stored under his or her name with date.

external_urls: stores URLs mentioned in replies or in the original post, temporarily disabled in the current crawler version.

Key functions:

start_crawling(): initializes MongoDB connections, extracts post URLs from the given forum webpage URL.

store_crawled_data(): by calling customized helper classes(refiner&parser,) three versions of the crawled data are stored on MongoDB server. Raw: stores the original HTML documents. Refined: stores parsed texts. Pure: stores tokenized texts without punctuation.

Helper classes:

refiner: handles the original HTML document and extracts post related information including date, username, original post contents, reply contents, etc. Each crawler has its own customized refiner class.

parser: uses NLTK and HTMLparser to tokenize post contents.

Key variables:

self.crawl_interval: this variable controls how long one crawler process will sleep after extracting web contents from one URL.

pageNum(under the main function): this variable controls the page range the crawler is going to crawl.

client(under the start_crawling function): this variable decides which database will be used to store the crawled datasets.

Library Versions:

nltk 3.6.1

pymongo 3.11.4

urllib3 1.25.8

beautifulsoup4 4.9.3

requests 2.25.1

Configure MongoDB server's remote connection: (MongoDB server version: 4.2.15)

sudo vim /etc/mongod.conf

```
frank@qifanz-MS-7C94:~$ sudo vim /etc/mongod.conf
```

change bindIp to 0.0.0.0; enable security authorization:

```
# network interfaces
net:
  port: 27017
  bindIp: 0.0.0.0

# how the process runs
processManagement:
  timeZoneInfo: /usr/share/zoneinfo

security:
  authorization: 'enabled'
```

To close remote connection:

change bindIp to 127.0.0.1; disable security authorization

```
# network interfaces
net:
  port: 27017
  bindIp: 127.0.0.1

# how the process runs
processManagement:
  timeZoneInfo: /usr/share/zoneinfo

#security:
# authorization: 'enabled'
```

Then restart Mongoddb server:

sudo service mongod restart

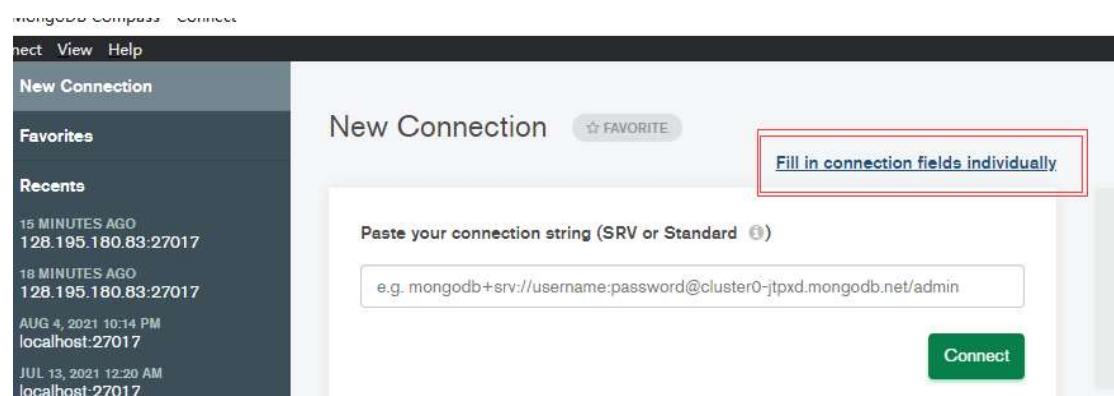
```
frank@qifanz-MS-7C94:~$ sudo service mongod restart
```

Admin: username: "admin", pwd: "ucidsplab_dbadmin"

Read-only: username: "db_viewer", pwd: "ucidsplab_dbviewer"

Remote connection by using MongoDB Compass (first time only):

Create a new connection, and select "Fill in connection fields individually":



Connect by using the Read-only username & password provided above:

The image shows a 'New Connection' dialog box with the following fields and values:

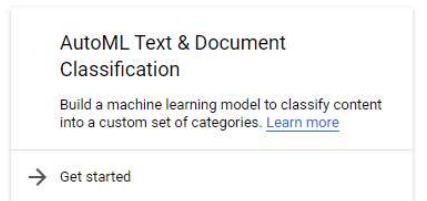
- Hostname:** 128.195.180.83
- Port:** 27017
- SRV Record:** ☐
- Authentication:** Username / Password
- Username:** db_viewer
- Password:** [Redacted]
- Authentication Database:** admin

A green 'Connect' button is located at the bottom right of the dialog box.

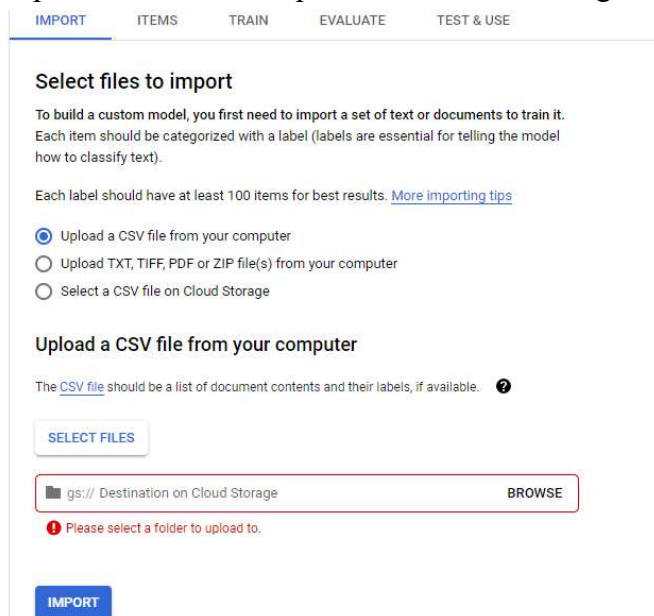
Google AutoML Natural Language:
Email Address: uci.dsp.dns.forum@gmail.com
Password: UCI_DSP_DNS_Forum

Select Text & Document Classification

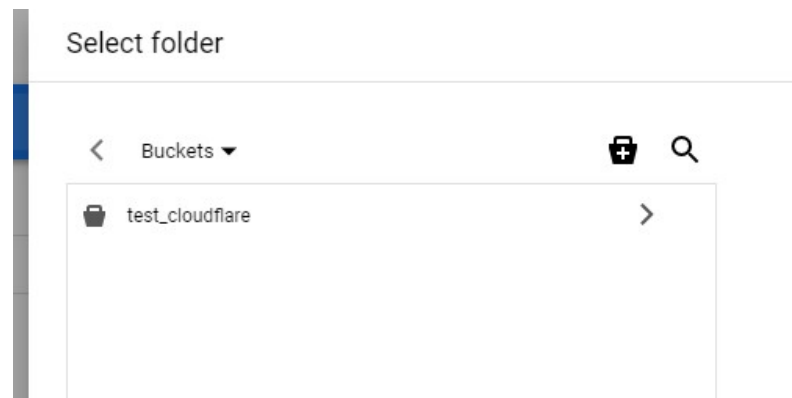
Natural Language products



Import CSV files and upload CSV files to Google Cloud Storage:



Choose an existing Cloud folder or create a new one:



After imports, assign labels to each data row and start training. The training process can take hours to complete.