

General Guidance

This is a guide with the information you need to build the solution. Here we present you the workflow we intent to execute in the web scraping and the expected outcome

Name: Totalexpress

Host: http://tracking.totalexpress.com.br/tracking/0?cpf_cnpj

Workflow

1- Access http://tracking.totalexpress.com.br/tracking/0?cpf_cnpj

2- Insert the given input in the “CPF”[“cpf”], “Primeiro Nome”[“name”] and “Cep”[“cep”] fields.
The input will be a python dictionary containing the “cpf”, “cep” and “name” keys



A screenshot of a web form for tracking a package. The form includes fields for "Primeiro Nome / Razão Social" (First Name / Business Reason), "CPF / CNPJ" (CPF/CNPJ), "Cep" (ZIP code), and "Número Verificador" (Verifier Number). A CAPTCHA challenge is displayed as a grid of numbers: 3, 8, 1, 7, 6. A red box highlights the "CPF / CNPJ" field, and the text "CPF Goes here" is written next to it. A yellow "pesquisar" (search) button is at the bottom right.

3- Crack the captcha challenge and insert the result in the “Número Verificador” field



A screenshot of a web form from TOTAL express. The form fields include "Primeiro Nome / Razão Social" (First Name / Business Reason), "CPF / CNPJ", "Cep", and "Número Verificador" (Verifier Number) which is set to "38176". Below these is a CAPTCHA field containing "38176" with a dashed border, accompanied by the text "Crack this captcha". A red arrow points from the text to the CAPTCHA field. At the bottom is a "pesquisar" (search) button.

4- Click the “Pesquisar” button

5- Collect the resulting Data

Collecting the resulting data

Once you have successfully entered the captcha challenge, you may encounter 3 possible outcomes:

- **Not found**
- **Information Found**

For each type of result, you will produce a JSON result. Your code must **always** return a python dictionary.

Not Found

You may encounter the following screens when executing the scraping:



Primeiro Nome / Razão Social:

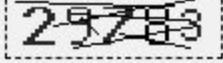
CPF / CNPJ:

Cep:

Número Verificador:
A dashed rectangular box containing the verification code "46a02".

Nenhuma encomenda encontrada! Tente novamente.



Primeiro Nome / Razão Social:	<input type="text"/>
CPF / CNPJ:	<input type="text"/>
Cep:	<input type="text"/>
Número Verificador:	<input type="text"/>
 25783	
<input type="button" value="pesquisar"/>	

CPF / CNPJ inválido!

If you received any of those screens, you must return a dictionary following the template:

```
{  
    "found_packages": False,  
    "total_packages": 0,  
    "packages": []  
}
```

Information Found

The expected result you will find when executing the scraper is the following:



Primeiro Nome / Razão Social:	<input type="text"/>
CPF / CNPJ:	<input type="text"/>
Cep:	<input type="text"/>
Número Verificador:	<input type="text"/> 4 891 5
<input type="button" value="pesquisar"/>	

Pedido	Nota Fiscal	AWB	Ver Detalhes
27690198012	15994	504527300	
25173009	736	104902440	
6280953601	559942	1826214835	

Table 1

When you reach that screen, you must click in all the rows in the “table 1” section and collect the information that will open up.

The screenshot shows a Google Chrome window for the Total Express Tracking website. At the top, it displays the URL: 'Nao seguro | tracking.totalexpress.com.br/tracking_encomenda.php?code=...'. The page features the Total Express logo and a red banner with the text 'Nossa entrega é Total!'. Below the banner is a table showing tracking history:

Data	Hora	Status
21/07/2016	15:29:22	ENCAMINHADO PARA O CENTRO DE DISTRIBUIÇÃO
21/07/2016	15:43:26	EM PROCESSO DE COLETA
22/07/2016	09:22:38	RECEBIDO NO CENTRO DE DISTRIBUIÇÃO
22/07/2016	14:26:45	TRANSFERENCIA PARA Recife/PE
25/07/2016	15:45:39	RECEBIDO NO CENTRO DE DISTRIBUIÇÃO Recife/PE
26/07/2016	10:04:46	SEPARADO PARA O ROTEIRO DE ENTREGA
26/07/2016	15:06:30	PROCESSO DE ENTREGA
26/07/2016	16:00:00	ENTREGA REALIZADA

To the right of the table is a search form with fields for 'Primeiro Nome / Razão Social', 'CPF / CNPJ', 'Cep', and 'Número Verificador' (with a placeholder '48425'). Below the form is a button labeled 'pesquisar'. Further down, there is a table with columns 'Pedido', 'Nota Fiscal', 'AWB', and 'Ver Detalhes' containing the values: 27690198012, 15994, 504527300; 25173009, 736, 104902440; and 6280953601, 559942, 1826214835.

This is the screen that will open up

You need to collect all the data and structure it in the JSON format. Each line in the “table 1” table will produce an item in the awway “packages”. (See the resulting JSON)

Data	Hora	Status
21/07/2016	15:29:22	ENCAMINHADO PARA O CENTRO DE DISTRIBUIÇÃO
21/07/2016	15:43:26	EM PROCESSO DE COLETA
22/07/2016	09:22:38	RECEBIDO NO CENTRO DE DISTRIBUIÇÃO
22/07/2016	14:26:45	TRANSFERENCIA PARA Recife/PE
25/07/2016	15:45:39	RECEBIDO NO CENTRO DE DISTRIBUIÇÃO Recife/PE
26/07/2016	10:04:46	SEPARADO PARA O ROTEIRO DE ENTREGA
26/07/2016	15:06:30	PROCESSO DE ENTREGA
26/07/2016	16:00:00	ENTREGA REALIZADA

For each line in the “Data” and “Status” columns, you will create a new entry in the `status_list` array in the resulting JSON.

If the “Status” line has the string “ENTREGA REALIZADA”, this is the “`delivery_date`” value.

Resulting JSON:

```
{
  "found_packages": true,
  "total_packages": 3,
```

```

"packages": [
    {
        "delivery_date": "26/07/2016",
        "package_id": "6280953601",
        "status_list": [
            {
                "date": "21/07/2016",
                "status": "ENCAMINHADO PARA O CENTRO DE DISTRIBUIÇÃO"
            },
            {
                "date": "21/07/2016",
                "status": "EM PROCESSO DE COLETA"
            },
            {
                "date": "22/07/2016",
                "status": "RECEBIDO NO CENTRO DE DISTRIBUIÇÃO"
            },
            {
                "date": "22/07/2016",
                "status": "TRANSFERENCIA PARA Recife/PE"
            },
            {
                "date": "25/07/2016",
                "status": "RECEBIDO NO CENTRO DE DISTRIBUIÇÃO Recife/PE"
            },
            {
                "date": "26/07/2016",
                "status": "SEPARADO PARA O ROTEIRO DE ENTREGA"
            },
            {
                "date": "26/07/2016",
                "status": "PROCESSO DE ENTREGA"
            },
            {
                "date": "26/07/2016",
                "status": "ENTREGA REALIZADA"
            }
        ]
    }
]
}

```

NOTE:

For simplicity sake, we omitted the other lines from the resulting JSON. But keep in mind that you need to scrap the data from **all the lines** in the “Table 1” table and add them to the “packages” array. We have attached a “result_example” containing a full example of accepted result

Implementation Details

You will receive a zip file containing files that will help you build the scraper code in the desired patterns. You need to fill the PES014.py file with your code. The file basically contains a skeleton of what we expect from a scraper that will be executed in our architecture. You can notice that there is a “test_request” method implemented. This method contains a base case of success and it may serve you as a guide to understand the final result we need your code to return.

For further details about how to build an acceptable scraper, access:

https://hydra.neurolake.io/hydra_sdk/index.html

https://hydra.neurolake.io/hydra_sdk/usage/basic_concepts.html

https://hydra.neurolake.io/hydra_sdk/usage/quickstart.html