

## ✓ 7.2 Data Collection through API

---

**Name:** Jann Moises Nyll B. De los Reyes

**Section:** CPE22S3

---

### Collecting temperature data from an API

#### About the data

In this notebook, we will be collecting daily temperature data from the [National Centers for Environmental Information \(NCEI\) API](#). We will use the Global Historical Climatology Network - Daily (GHCND) data set; see the documentation [here](#).

## ✓ Using the NCEI API

Paste your token below.

```
1 import requests
2
3 def make_requests(endpoint,payload=None):
4     """
5     Make a request to a specific endpoint on the weather API
6     passing headers and optional payload.
7
8     Parameters:
9         - endpoint: The endpoint of the API you want to
10             make a GET request to.
11         - payload: A dictionary of data to pass along
12             with the request.
13
14     Returns:
15         Response object.
16     """
17
18     return requests.get(f'https://www.ncdc.noaa.gov/cdo-web/api/v2/{endpoint}',
19                         headers={
20                             'token': 'fbxFRxUcYeGBgSMcECTkBoAb1KBwVBqd'
21                         },
22                         params = payload
23 )
```

## ✓ See what datasets are available

We can make requests to the datasets endpoint to see what datasets are available. We also pass in a dictionary for the payload to get datasets that have data after the start date of October 1, 2018.

```
1 # see what dataset are available
2 response = make_requests('datasets',{'startdate':'2018-10-01'})
3 response.status_code

200

1
2 response = make_requests('datasets',{'startdate':'2024-01-01'})
3 response.status_code

200
```

Status code of 200 means everything is OK. More codes can be found [here](#).

## ✓ Get the keys of the result

The result is a JSON object which we can access with the `json()` method of our `Response` object. JSON objects can be treated like dictionaries, so we can access the `keys()` just like we would a dictionary:

```
1 response.json().keys()

dict_keys(['metadata', 'results'])
```

The `metadata` of the JSON response will tell us information about the request and data we got back:

```
1 response.json()['metadata']

{'resultset': {'offset': 1, 'count': 11, 'limit': 25}}
```

## ✓ Figure out what data is in the result

The `results` key contains the data we requested. This is a list of what would be rows in our dataframe. Each entry in the list is a dictionary, so we can look at the keys to get the fields:

```
1 response.json()['results'][10]

{'uid': 'gov.noaa.ncdc:C00313',
 'mindate': '1900-01-01',
 'maxdate': '2014-01-01',
 'name': 'Precipitation Hourly',
 'datacoverage': 1,
 'id': 'PRECIP_HLY'}

1 response.json()['results'][0].keys()

dict_keys(['uid', 'mindate', 'maxdate', 'name', 'datacoverage', 'id'])
```

## ✓ Parse the result

We don't want all those fields, so we will use a list comprehension to take only the `id` and `name` fields out

```
1 [(data['id'], data['name']) for data in response.json()['results']]

[('GHCND', 'Daily Summaries'),
 ('GSOM', 'Global Summary of the Month'),
 ('GSOY', 'Global Summary of the Year'),
 ('NEXRAD2', 'Weather Radar (Level II)'),
 ('NEXRAD3', 'Weather Radar (Level III)'),
 ('NORMAL_ANN', 'Normals Annual/Seasonal'),
 ('NORMAL_DLY', 'Normals Daily'),
 ('NORMAL_HLY', 'Normals Hourly'),
 ('NORMAL_MLY', 'Normals Monthly'),
 ('PRECIP_15', 'Precipitation 15 Minute'),
 ('PRECIP_HLY', 'Precipitation Hourly')]
```

## ✓ Figure out which data category we want

The `GHCND` data containing daily summaries is what we want. Now we need to make another request to figure out which data categories we want to collect. This is the `datacategories` endpoint. We have to pass the `datasetid` for `GHCND` as the payload so the API knows which dataset we are asking about

```
1 #get data category id
2 response = make_requests(
3     'datacategories',
4     payload={
5         'datasetid': 'GHCND'
6     }
7 )
8 response.status_code

200
```

Since we know the API gives us a `metadata` and a `results` key in each response, we can see what is in the `results` portion of the JSON response:

```
1 response.json()['results']

[{'name': 'Evaporation', 'id': 'EVAP'},
 {'name': 'Land', 'id': 'LAND'},
 {'name': 'Precipitation', 'id': 'PRCP'},
 {'name': 'Sky cover & clouds', 'id': 'SKY'},
 {'name': 'Sunshine', 'id': 'SUN'},
 {'name': 'Air Temperature', 'id': 'TEMP'},
 {'name': 'Water', 'id': 'WATER'},
 {'name': 'Wind', 'id': 'WIND'},
 {'name': 'Weather Type', 'id': 'WXTYPE'}]
```

## ✓ Grab the data type ID for the Temperature category

We will be working with temperatures, so we want the `TEMP` data category. Now, we need to find the `datatypes` to collect. For this, we use the `datatypes` endpoint and provide the `datacategoryid` which was `TEMP`. We also specify a limit for the number of `datatypes` to return with the payload. If there are more than this we can make another request later, but for now, we just want to pick a few out:

```
1 response = make_requests(
2     'datatypes',
3     payload={
4         'datacategoryid': 'TEMP',
5         'limit': 100
6     }
7 )
8 response.status_code

200
```

We can grab the `id` and `name` fields for each of the entries in the `results` portion of the data. The fields we are interested in are at the bottom:

```
1 [(datatype['id'], datatype['name']) for datatype in response.json()['results']][-5:] #look at the last 5
[('MNTM', 'Monthly mean temperature'),
 ('TAVG', 'Average Temperature.'),
 ('TMAX', 'Maximum temperature'),
 ('TMIN', 'Minimum temperature'),
 ('TOBS', 'Temperature at the time of observation')]
```

Now that we know which `datatypes` we will be collecting, we need to find the location to use. First, we need to figure out the location category.

This is obtained from the `locationcategories` endpoint by passing the `datasetid` :

```
1 #get location category id
2 response = make_requests(
3     'locationcategories',
4     {
5         'datasetid': 'GHCND'
6     }
7 )
8 response.status_code

200
```

We can use `pprint` to print dictionaries in an easier-to-read format. After doing so, we can see there are 12 different location categories, but we are only interested in `CITY` :

```
1 import pprint
2 pprint.pprint(response.json())

{'metadata': {'resultset': {'count': 12, 'limit': 25, 'offset': 1}},
 'results': [{'id': 'CITY', 'name': 'City'},
              {'id': 'CLIM_DIV', 'name': 'Climate Division'},
              {'id': 'CLIM_REG', 'name': 'Climate Region'},
              {'id': 'CNTRY', 'name': 'Country'},
              {'id': 'CNTY', 'name': 'County'},
              {'id': 'HYD_ACC', 'name': 'Hydrologic Accounting Unit'},
              {'id': 'HYD_CAT', 'name': 'Hydrologic Cataloging Unit'},
              {'id': 'HYD_REG', 'name': 'Hydrologic Region'},
              {'id': 'HYD_SUB', 'name': 'Hydrologic Subregion'},
              {'id': 'ST', 'name': 'State'},
              {'id': 'US_TERR', 'name': 'US Territory'},
              {'id': 'ZIP', 'name': 'Zip Code'}]}
```

## ✓ Get NYC Location ID

In order to find the location ID for New York, we need to search through all the cities available. Since we can ask the API to return the cities sorted, we can use binary search to find New York quickly without having to make many requests or request lots of data at once. The following function makes the first request to see how big the list of cities is and looks at the first value. From there it decides if it needs to move towards the beginning or end of the list by comparing the city we are looking for to others alphabetically. Each time it makes a request it can rule out half of the remaining data to search.

```
1
```

```

1
2 def get_item(name, what, endpoint, start=1, end=None):
3     """
4     Grab the JSON payload for a given field by name using binary search.
5
6     Parameters:
7         - name : data item to look for
8         - what : specification in the dictionary.Dictionary specifying what the item in `name` is.
9         - endpoint :Where to look for the item.
10        - start : beginning of the set. The position to start at. We don't need to touch this, but the
11        function will manipulate this with recursion.
12        - end : The last position of the cities. Used to find the midpoint, but
13        like `start` this is not something we need to worry about.
14
15    Returns:
16        - Dictionary of the information for the item if found otherwise
17        an empty dictionary.
18
19    """
20
21    # Find the midpoint of the dataset which we can use to cut the data in half each time
22    mid = (start + (end if end else 1)) // 2
23
24    # change name input to lowercase, so this is not case-sensitive
25    name = name.lower()
26
27    #define the payload we will send each request
28    payload = {
29        'datasetid':'GHCND',
30        'sortfield':'name',
31        'offset': mid, #we will change the offset each time
32        'limit' : 1 #we only want one value back
33    }
34
35    #make our request adding any optional filter parameter 'what'
36    response = make_requests(endpoint, **payload, **what})
37
38    if response.ok:
39        # Get the end value , if response is ok, grab the end index from the response metadata the first time through
40        end = end if end else response.json()['metadata']['resultset']['count']
41
42        #Grab the lowercase version of the current name
43        current_name = response.json()['results'][0]['name'].lower()
44
45        # if what we are searching for is in the current name, we have found our item
46        if name in current_name:
47            return response.json()['results'][0] # return the found item
48        else:
49            # if our start index is greater than or equal to our end, we couldn't find it
50            if start >= end:
51                return {}
52            # our name comes before the current name in the alphabet, so we search further to the left
53            elif name < current_name:
54                return get_item(name, what, endpoint, start, mid - 1)
55            # our name comes after the current name in the alphabet, so we search further to the right
56            elif name > current_name:
57                return get_item(name, what, endpoint, mid+1, end)
58        else:
59            # response wasn't ok, use code to determine why
60            print(f'Response not ok, status {response.status_code}')
61
62 def get_location(name):
63     """
64     Grab the JSON payload for the location by name using binary search.
65     Parameters:
66         - name: The city to look for.
67     Returns:
68         Dictionary of the information for the city if found otherwise
69         an empty dictionary.
70     """
71     return get_item(name, {'locationcategoryid':'CITY'}, 'locations')

```

When we use binary search to find New York, we find it in just 8 requests despite it being close to the middle of 1,983 entries

```

1 # get NYC id
2 nyc = get_location('New York')
3 nyc

```

```

{'mindate': '1869-01-01',
 'maxdate': '2024-03-14',
 'name': 'New York, NY US',
 'datacoverage': 1,
 'id': 'CITY:US360019'}

```

## Get the station ID for Central Park

The most granular data is found at the station level:

```

1 central_park = get_item('NY City Central Park', {'locationid':nyc['id']}, 'stations')
2 central_park

```

```
{'elevation': 42.7,
'mindate': '1869-01-01',
'maxdate': '2024-03-13',
'latitude': 40.77898,
'name': 'NY CITY CENTRAL PARK, NY US',
'datacoverage': 1,
'id': 'GHCND:USW00094728',
'elevationUnit': 'METERS',
'longitude': -73.96925}
```

## Request the temperature data

Finally, we have everything we need to make our request for the New York temperature data. For this we use the `data` endpoint and provide all the parameters we picked up throughout our exploration of the API:

```
1 # get NYC daily summaries data
2 response = make_requests(
3     'data',
4     {
5         'datasetid' : 'GHCND',
6         'stationid' : central_park['id'],
7         'locationid' : nyc['id'],
8         'startdate' : '2018-10-01',
9         'enddate' : '2018-10-31',
10        'datatypeid' : ['TMIN','TMAX','TOBS'], #temperature at time of observation,min, and max
11        'units' : 'metric',
12        'limit' : 1000
13    }
14 )
15 response.status_code

200
```

## Create a DataFrame

The Central Park station only has the daily minimum and maximum temperatures.

```
1 import pandas as pd
2
3 df = pd.DataFrame(response.json()['results'])
4 df.head()
```

	date	datatype	station	attributes	value
0	2018-10-01T00:00:00	TMAX	GHCND:USW00094728	„W,2400	24.4
1	2018-10-01T00:00:00	TMIN	GHCND:USW00094728	„W,2400	17.2
2	2018-10-02T00:00:00	TMAX	GHCND:USW00094728	„W,2400	25.0
3	2018-10-02T00:00:00	TMIN	GHCND:USW00094728	„W,2400	18.3
4	2018-10-03T00:00:00	TMAX	GHCND:USW00094728	„W,2400	23.3

We didn't get `TOBS` because the station doesn't measure that:

```
1 df.datatype.unique()

array(['TMAX', 'TMIN'], dtype=object)
```

Despite showing up in the data as measuring it... Real-world data is dirty!

```
1 if get_item(
2     'NY City Central Park',{'locationid' : nyc['id'], 'datatypeid' : 'TOBS'}, 'stations'
3 ):
4 ):
5
6     print('Found!')

Found!
```

## Using a different station

Let's use LaGuardia airport instead. It contains `TAVG` (average daily temperature):

```
1 laguardia = get_item(
2     'LaGuardia', {'locationid' : nyc['id']], 'stations'
3 )
4
5 laguardia

{'elevation': 3,
'mindate': '1939-10-07',
'maxdate': '2024-03-14',
'latitude': 40.77945,
'name': 'LAGUARDIA AIRPORT, NY US',
```

```
'datacoverage': 1,
'id': 'GHCND:USW00014732',
'elevationUnit': 'METERS',
'longitude': -73.88027}
```

We make our request using the LaGuardia airport station this time and ask for TAVG instead of TOBS .

```
1 #get NYC daily summaries data
2 response = make_requests(
3     'data',
4     {
5         'datasetid' : 'GHCND',
6         'stationid' : laguardia['id'],
7         'location' : nyc['id'],
8         'startdate' : '2018-10-01',
9         'enddate' : '2018-10-31',
10        'datatypeid' : ['TMIN', 'TMAX', 'TAVG'], # temperature at time observation, min, and max
11        'units' : 'metric',
12        'limit' : 1000
13    }
14 )
15 )
16
17 response.status_code

200
```

The request was successful, so let's make a dataframe:

```
1 df = pd.DataFrame(response.json()['results'])
2 df.head()
3
```

	date	datatype	station	attributes	value
0	2018-10-01T00:00:00	TAVG	GHCND:USW00014732	H,S,	21.2
1	2018-10-01T00:00:00	TMAX	GHCND:USW00014732	,W,2400	25.6
2	2018-10-01T00:00:00	TMIN	GHCND:USW00014732	,W,2400	18.3
3	2018-10-02T00:00:00	TAVG	GHCND:USW00014732	H,S,	22.7
4	2018-10-02T00:00:00	TMAX	GHCND:USW00014732	,W,2400	26.1

We should check we got what we wanted: 31 entries for TAVG, TMAX and TMIN (1 per day):

```
1 df.datatype.value_counts()

datatype    31
TAVG        31
TMAX        31
TMIN        31
Name: datatype, dtype: int64
```

Write the data to a CSV file for use in other notebooks.

```
1 df.to_csv('datos/nyc_temperatures.csv', index = False)
```