# Handout_Data Collection III

This week, we will apply what we have learned to another website, Google review. With this practice, we are going to collect tourists' reviews about a destination on Google review using webdriver. This approach allows us to collect data from dynamic webpage.

1. Create a new Python file as you did in Data Collection I and II. Please check the site below for your reference.
   https://chunshengj.github.io/579-class/data_analysis/Install_Anaconda_Jupyter_Package.html

2. Import packages and set working directory

```python
import os
import requests
import bs4
import pandas as pd
from selenium import webdriver
from selenium.webdriver.chrome.service import Service
from selenium.webdriver.chrome.options import Options
from selenium.webdriver.common.by import By
from selenium.webdriver.support.ui import WebDriverWait
from selenium.webdriver.support import expected_conditions as EC
import time
```

```python
os.chdir("/Users/sheng/Jupyter/HON322M")  # for Mac
# os.chdir("C:\\Users\\sheng\\Jupyter\\HON322M")  # for Windows
```

3. Select the webpage we collect and scrape the content in the webpage
   1) Go to the website below
      a. https://www.google.com/maps/place/Elephant+Nature+Park/@19.2169619,98.8598287,18z/data=!4m10!1m2!2m1!1selephant+nature+park!3m6!1s0x30da3aa86b496d8f:0xdb2c200fd3b02b15!8m2!3d19.2164372!4d98.8613856!9m1!1b1

b. The specific information we want to collect includes 1) Reviewer name, 2) Review rating, 3) Review date, and 4) Review

2) Assign the URL to a variable name "url (copy and paste the link provided in 3.1)

```
url="https://www.google.com/maps/place/Elephant+Nature+Park/@19.2169619,98.8598287,18z/data=!4m10!1m2!2m1!1selephant
```

3) Open the webpage to collect review information with the "selenium" package

    a. Inspect to identify the elements that contain the review information



    b. Incomplete review before clicking 'More'



Need to click this "More" to view the complete review

```
▼<div class="GHT2ce"> flex
  ▶<div class="DU9Pgb"> ⋯ </div> flex
  ▼<div>
    ▼<div class="MyEned" tabindex="-1" id="ChZDSUhNMG9nS0V
      JQ0FnSUQ5NXYtM2VnEAE" lang="en">
      ▼<span class="wiI7pd"> == $0
          "One of the few really ethical elephant rescue
          parks. Our guide En was amazing, she seemed to
          enjoy every moment just like us. I can only
          recommend it. I know the price is high, but they
          have a lot of …"
```
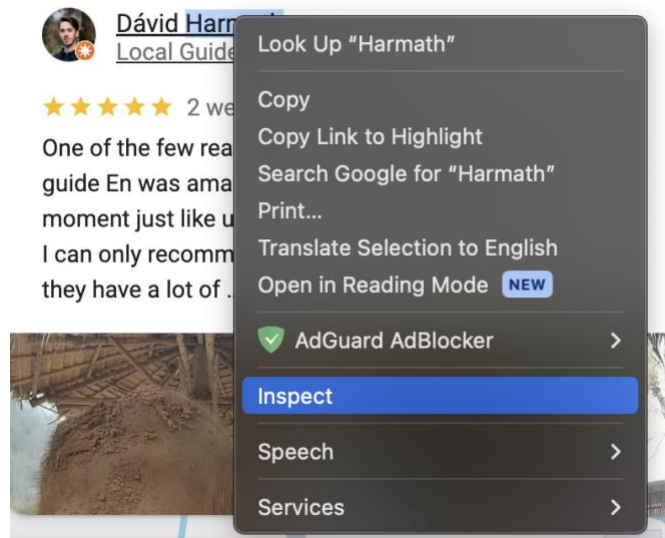
c.  After clicking "More", the review is complete.



Dávid Harmath
Local Guide · 223 reviews · 757 photos

★★★★★  2 weeks ago  NEW

One of the few really ethical elephant rescue parks. Our guide En was amazing, she seemed to enjoy every moment just like us.
I can only recommend it. I know the price is high, but they have a lot of expenses because elephants eat a lot.
The buffet lunch was quite good too, there were a lot of options and everything was vegan.



```
▼<div class="GHT2ce"> flex
  ▶<div class="DU9Pgb"> ⋯ </div> flex
  ▼<div>
    ▼<div class="MyEned" tabindex="-1" id="ChZDSUhNMG9nS0V
      JQ0FnSUQ5NXYtM2VnEAE" lang="en">
      ▼<span class="wiI7pd"> == $0
          "One of the few really ethical elephant rescue
          parks. Our guide En was amazing, she seemed to
          enjoy every moment just like us. I can only
          recommend it. I know the price is high, but they
          have a lot of expenses because elephants eat a
          lot. The buffet lunch was quite good too, there
          were a lot of options and everything was vegan."
```

d.  In this case, "request" commands are not usable to collect a complete review. So, we need to obtain the html content of the webpage using webdriver from "selenium" package.

e.  Install the Selenium package
    i. See the section 1 above for package installation

ii. Make sure that you have the selenium package installed

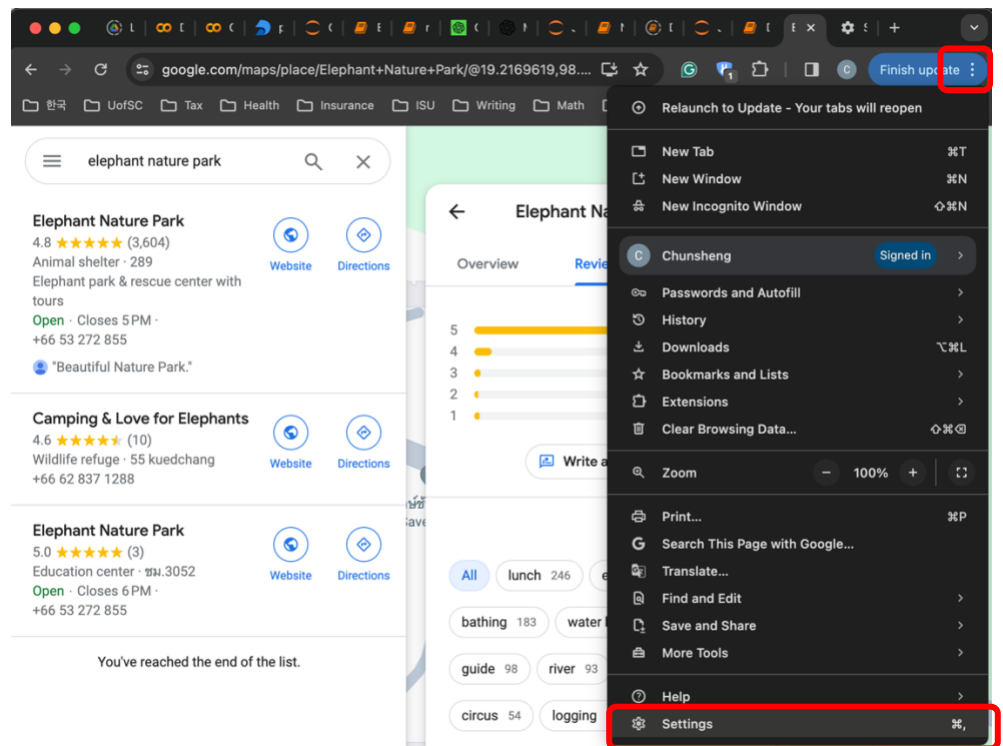f. Import different selenium functions in your python file

```python
import os
import requests
import bs4
import pandas as pd
from selenium import webdriver
from selenium.webdriver.chrome.service import Service
from selenium.webdriver.chrome.options import Options
from selenium.webdriver.common.by import By
from selenium.webdriver.support.ui import WebDriverWait
from selenium.webdriver.support import expected_conditions as EC
import time
```

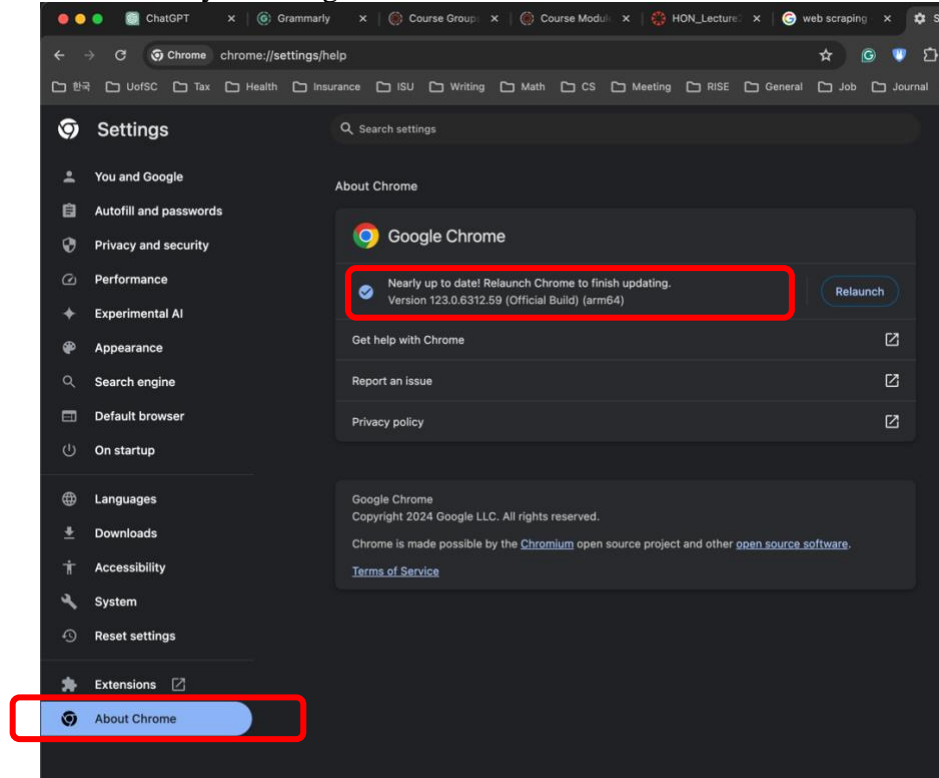g. Add a couple of command lines for selenium options after "url" defined previously

```python
url="https://www.google.com/maps/place/Elephant+Nature+Park/@19.2169619,98.8598287,18z/data=!4m10!1m2!2m1!1selephant
```

```python
options = webdriver.ChromeOptions()
options.add_experimental_option("detach", True)
```

h. Then, we need to open the webpage via selenium to make the contents of the webpage collectable. We will open the webpage to collect with Chrome. To do that, we need "Chromedriver"

i. Check your version of Chrome

i. Open Chrome and click three dots at the end of the right side of the search box. Then, hit "Settings"

ii. Click "About Chrome" in the bottom of the menu to the right
iii. Check your Google Chrome version



j. Go to https://chromedriver.chromium.org/downloads to download the chromedriver for your version. My version is 123.

**Stable**

Version: 123.0.6312.86 (r1262506)

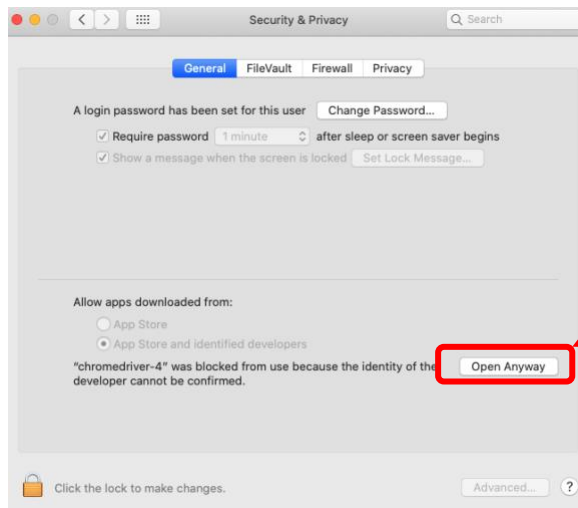| Binary | Platform | URL | HTTP status |
|---|---|---|---|
| chrome | linux64 | https://storage.googleapis.com/chrome-for-testing-public/123.0.6312.86/linux64/chrome-linux64.zip | 200 |
| chrome | mac-arm64 | https://storage.googleapis.com/chrome-for-testing-public/123.0.6312.86/mac-arm64/chrome-mac-arm64.zip | 200 |
| chrome | mac-x64 | https://storage.googleapis.com/chrome-for-testing-public/123.0.6312.86/mac-x64/chrome-mac-x64.zip | 200 |
| chrome | win32 | https://storage.googleapis.com/chrome-for-testing-public/123.0.6312.86/win32/chrome-win32.zip | 200 |
| chrome | win64 | https://storage.googleapis.com/chrome-for-testing-public/123.0.6312.86/win64/chrome-win64.zip | 200 |
| chromedriver | linux64 | https://storage.googleapis.com/chrome-for-testing-public/123.0.6312.86/linux64/chromedriver-linux64.zip | 200 |
| chromedriver | mac-arm64 | https://storage.googleapis.com/chrome-for-testing-public/123.0.6312.86/mac-arm64/chromedriver-mac-arm64.zip | 200 |
| chromedriver | mac-x64 | https://storage.googleapis.com/chrome-for-testing-public/123.0.6312.86/mac-x64/chromedriver-mac-x64.zip | 200 |
| chromedriver | win32 | https://storage.googleapis.com/chrome-for-testing-public/123.0.6312.86/win32/chromedriver-win32.zip | 200 |
| chromedriver | win64 | https://storage.googleapis.com/chrome-for-testing-public/123.0.6312.86/win64/chromedriver-win64.zip | 200 |

    k.  Install Chromedriver (PC)
- iv. Copy and paste the link on a new window to download chromedriver (win32 or 64.zip depending on your PC)
- v. Unzip the zipped folder downloaded. Then, find chromedriver.exe
- vi. Move the exe file to the working directory for your python file

    l.  Install Chromedriver (Mac)
- vii. Download chromedriver (mac64.zip or mac-arm64.zip (for M1))
- viii. Double click the downloaded file

    m.  If you have the following message, you can move to the next step

```
Last login: Sat Mar 16 16:04:40 on ttys000
/Users/sheng/Jupyter/HON322M/chromedriver ; exit;
(base) sheng@gimchunseongdeMac-Studio ~ % /Users/sheng/Jupyter/HON322M/chromedri
ver ; exit;
Starting ChromeDriver 123.0.6312.58 (6b4b19e9dfbb93aa414dc045bd445287977d8d7a-re
fs/branch-heads/6312_46@{#3}) on port 9515
Only local connections are allowed.
Please see https://chromedriver.chromium.org/security-considerations for suggest
ions on keeping ChromeDriver safe.
ChromeDriver was started successfully.
```

n. If the file is not opened, so you have the following message,



You should go to "System Preference" => "Security & Privacy" => "General"



Hit "Open Anyway" to make the driver usable

o. Open the blank Chrome and then the webpage to collect
p. There are two ways to open the Chrome. Try both and find the one that works on your computer.

ix. Use executable_path for the first argument of webdriver.Chrome

```
driver = webdriver.Chrome(executable_path = "/Users/sheng/Jupyter/HON322M/chromedriver",
                          options = options)
driver.get(url)
```

x. Use Service function

```
s = Service("/Users/sheng/Jupyter/HON322M/chromedriver")
driver = webdriver.Chrome(service = s, options = options)

driver.get(url)
```

xi. New google chrome window will be opened after executing "drver.get(url)"

4. Collect review information
   1) Click all "More" to have complete review information

a. Inspect element





b. Identify a code pattern
   i. Element: button
   ii. Attribute: class
   iii. Attribute value: "w8nwRe kyuRq" => This value can be changed. You should check the attribute value by yourself.
   iv. This 'button' element contains text "More".

c. Use find_elements(By.XPATH) to identify the elements and check the number of "More" needed to click

```
viewMore=driver.find_elements(By.XPATH, "//button[./text()='More']")
len(viewMore)
```

7

d. Locate each "More" and click it

```
for i in range(0,len(viewMore)):
    invisible = WebDriverWait(driver, 15).until(
        EC.presence_of_element_located((By.XPATH, "//button[./text()='More']")))
    driver.execute_script("arguments[0].click();", invisible)
    time.sleep(2)
```

C  Conor Crews
   34 reviews · 9 photos

★★★★★  a month ago

Had a great experience walking around the park,
meeting all the elephants and watching them
take a bath. We did the overnight tour, so that
was our first day along with an outstanding
buffet for lunch and dinner.

The morning of the second day, we made some
food for the elephants, took a short walk outside
of the park to see the other elephants, and then
made an elephant cake to feed to one of the
older elephants. We also had an amazing
breakfast and lunch buffet again!

Highly recommended this place as it was ethical
and all the animals are treated with respect and
cared for. No riding or bathing them was allowed.
All animals are rescued and brought back to
good health.

Now, the reviews are fully expanded, and we can start collect each information.

2) Assign the HTML document to a variable "bs_obj"

```
bs_obj = bs4.BeautifulSoup(driver.page_source, "html.parser")
```

3) Identify the review information
   a. Inspect element

b. We can see three different contents in this tag (element: div, attribute: class, and attribute value: "jJc9Ad"), and all the information we need is stored under this tag. As mentioned before, the attribute value can be changed, so we should always check it.

c. Find all tags with element: div, attribute: class, and attribute value: "jJc9Ad" and assign them to a variable "reviews" using the findAll function.

```
reviews = bs_obj.findAll("div", {"class": "jJc9Ad"})
print(len(reviews))

8
```

There are 8 reviews on the page (if we scroll down, more reviews can be loaded and displayed).

4) Reviewer
   a. Inspect element

b. We can use find function to find div element, class attribute, and "d4r55" attribute value to get reviewer name

```
# reviewer name
reviews[0].find("div", {"class": "d4r55"}).text
```

```
'Dávid Harmath'
```

5) Review date
   a. Inspect element
   b. We can see the review date information is under span element, class attribute, and "rsqaWe" attribute value



   c. Similarly, we can use find function to extract review date

```
# review_date
reviews[0].find("span", {"class": "rsqaWe"}).text
```

```
'2 weeks ago'
```

6) Review rating
   a. Inspect element



   b. The review rating is 5 and this information is the attribute value, so we need to access 'aria-label' attribute to get its value. Let's see how we can get this information step by step.
      i. Find span element

```
# review rating
reviews[0].find("span", {"class": "kvMYJc"})
```

```
<span aria-label="5 stars" class="kvMYJc" role="img"><span class="hCCjke vzX5Ic google-symbols NhBTye">􀀀</span><spa
n class="hCCjke vzX5Ic google-symbols NhBTye">􀀀</span><span class="hCCjke vzX5Ic google-symbols NhBTye">􀀀</span><sp
an class="hCCjke vzX5Ic google-symbols NhBTye">􀀀</span><span class="hCCjke vzX5Ic google-symbols NhBTye">􀀀</span></
span>
```

    ii.  Get "aria-label" attribute

```
reviews[0].find("span", {"class": "kvMYJc"})['aria-label']
```

```
'5 stars'
```

    iii.  We can use split function to collect only the numeric value 5.

```
reviews[0].find("span", {"class": "kvMYJc"})['aria-label'].split(' ')
```

```
['5', 'stars']
```

    iv.  Now, it's split into two strings in a list.
    v.  Then, we just need to select the first one based on the index

```
reviews[0].find("span", {"class": "kvMYJc"})['aria-label'].split(' ')[0]
```

```
'5'
```

    vi.  The value is enclosed in single quotes, indicating that it is a string type rather than a numeric type. To convert the string to a numeric type, we can use int function. So, the final code to collect review rating is:

```
# review rating
int(reviews[0].find("span", {"class": "kvMYJc"})['aria-label'].split(' ')[0])
```

```
5
```

7) Review
   a.  Inspect element and find review



```
# review
reviews[0].find("span", {"class": "wiI7pd"}).text
```

```
'One of the few really ethical elephant rescue parks. Our guide En was amazing, she seemed to enjoy every moment ju
st like us.\nI can only recommend it. I know the price is high, but they have a lot of expenses because elephants e
at a lot.\nThe buffet lunch was quite good too, there were a lot of options and everything was vegan.'
```

8) For now, we have collected all the review information for a single review. As there are eight reviews in this webpage (scrolling down to view more reviews), we need to collect all the information of the eight reviews with the for loop.

   a. When we run a for loop, the review information in variables will be overwritten whenever we collect the information from different reviews. As a result, the final return shows the information of the last review. So, we should store the information of each review in lists (containers) for every round of the loop.

      i. Create empty lists for variables
      ii. Under for loop, create four variables (here, I named the four variables as reviewer, review_date, review_rating, and text) and assign the values (the code above used to collect the review information) to the four variables
      iii. Use "append" function to add the four variables to the empty lists
      iv. Print each variable to view the result

**i.**
```python
reviewers, review_dates, review_ratings, texts = [], [], [], []
```

**ii.**
```python
for i in reviews:

    reviewer = i.find("div", {"class": "d4r55"}).text
    review_date = i.find("span", {"class": "rsqaWe"}).text
    review_rating = int(i.find("span", {"class": "kvMYJc"})['aria-label'].split(' ')[0])
    text = i.find("span", {"class": "wiI7pd"}).text.strip()
```

**iii.**
```python
    # append the values collected to the empty lists
    reviewers.append(reviewer)
    review_dates.append(review_date)
    review_ratings.append(review_rating)
    texts.append(text)
```

```python
    # pring the result if we want to check the values we collected
    print(reviewer)
    print(review_date)
    print(review_rating)
    print(text)
    print("--" * 50 + '\n')
```

**iv.**
```
Dávid Harmath
2 weeks ago
5
One of the few really ethical elephant rescue parks. Our guide En was amazing, she seemed to enjoy every moment jus
t like us.
I can only recommend it. I know the price is high, but they have a lot of expenses because elephants eat a lot.
The buffet lunch was quite good too, there were a lot of options and everything was vegan.
----------------------------------------------------------------------------------------------------

Diane Rogers
2 weeks ago
5
Highly recommend coming here — excellent facilities, well maintained and amazing care for their animals (elephants,
dogs, cats). I felt that the money paid went to great cause. We did the 1/2 day morning, so was anxious when I saw
crowds arriving by bus. However, they formed small groups of approx 10-12 people with their own guide, who did the
2 hour walk and education (MINT was excellent). Very enjoyable and food at lunch was terrific, better than any rest
aurant! Highly suggest you come here versus any other 'sanctuary' that still does elephant rides or feeding.
----------------------------------------------------------------------------------------------------
```

9) Create a dataframe (table with columns and rows) for the review data collected

```
# create a dataframe
df = pd.DataFrame({"reviewer": reviewers, "review_date": review_dates, "review_rating": review_ratings,
                   "text": texts})
df
```

| | reviewer | review_date | review_rating | text |
|---|---|---|---|---|
| 0 | Dávid Harmath | 2 weeks ago | 5 | One of the few really ethical elephant rescue ... |
| 1 | Diane Rogers | 2 weeks ago | 5 | Highly recommend coming here - excellent facil... |
| 2 | Jodie Bowpitt | 2 weeks ago | 5 | The most beautiful experience. Por our tour gu... |
| 3 | Izzy the Bricky | a month ago | 5 | The most amazing and real place to see elephan... |
| 4 | Sean Michael Yip | a week ago | 5 | This is one of the few truly ethical elephant ... |
| 5 | Benjamin Ryberg | 2 months ago | 5 | The most spectacular and unforgettable experie... |
| 6 | Ilse Schouten | a month ago | 5 | We did the overnight tour, and it was amazing.... |
| 7 | Conor Crews | a month ago | 5 | Had a great experience walking around the park... |

10) Export data to our local computer in a csv format

```
# export the dataframe to a csv file
df.to_csv("Google review.csv")
```

a. The exported file is stored in the folder that we set a working directory before
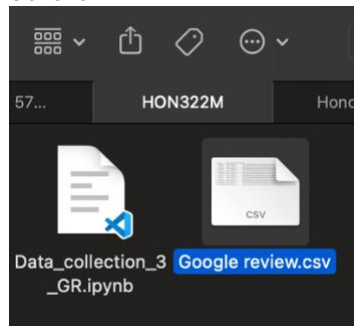


15

View　Zoom　Add Category　Pivot Table　Insert　Table　Chart　Text　Shape　Media　Comment　Share　Format　Organize

Sheet 1

⊘ Table data was imported and can be adjusted. ›

| | reviewer | Google | review_date | review_rating | text |
|---|---|---|---|---|---|
| 0 | Dávid Harmath | Google | 2 weeks ago | 5 | One of the few really ethical elephant rescue parks. Our guide En was amazing, she seemed to enjoy every moment just like us. I can only recommend it. I know the price is high, but they have a lot of expenses because elephants eat a lot. The buffet lunch was quite good too, there were a lot of options and everything was vegan. |
| 1 | Diane Rogers | Google | 2 weeks ago | 5 | Highly recommend coming here - excellent facilities, well maintained and amazing care for their animals (elephants, dogs, cats). I felt that the |
| 2 | Jodie Bowpitt | Google | 2 weeks ago | 5 | The most beautiful experience. Por our tour guide was amazing and so knowledgeable. A community that really looks after the elephants. We |
| 3 | Izzy the Bricky | Google | a month ago | 5 | The most amazing and real place to see elephants living in their natural habitat! Elephant Nature Park, is one of very few places in Thailand w Fantastic place, beautiful scenery, great guide and tasty food! |
| 4 | Sean Michael Yip | Google | a week ago | 5 | This is one of the few truly ethical elephant scantuaries. Lunch and hotel pick ups are provided.

There is no showering or feeding of the elephants. They are truly allowed to roam the park to their hearts content.

They are not forced to lie down and get showered by tourists. Any "scantuary" that does this has to train and force the elephants to endure e:

This park is beautifully constructed and you can see the true care given to the elephants. Many of which have serious disabilities from their pi

The food provided was vegan and delicious. Coffee and water are provided for free and the drinks/snacks you can purchase are very reasona |
| 5 | Benjamin Ryberg | Google | 2 months ago | 5 | The most spectacular and unforgettable experience of my two-week trip to Thailand (which was filled with stellar experiences). I did the full-d

Spent the morning with a small group, guides, and three magnificent elephants! My understanding is that this is one of the few experiences a

The guides (and photographer!) were fantastic and spoke excellent English. The included lunch was very tasty.

Exploring the actual park in the afternoon was wonderful as well, but getting to bond with the three elephants in the morning made the experi |
| 6 | Ilse Schouten | Google | a month ago | 5 | We did the overnight tour, and it was amazing. Started the first day around the main area, learning about the elephants, seeing them go for a

The breakfast/lunch/dinners were incredible. I'm pretty sure it was completely vegan, but with really good imitation chicken etc, so it didn't fee

Overall highly recommend doing overnight tour as I feel like you get some free time to just watch the elephants and walk around. And you see |

16