Entrée [ ]:      1

## Table of content

Entrée [ ]:      1

Entrée [1]:

```python
import pandas as pd
import numpy as np

import matplotlib.pyplot as plt

from slugify import slugify
import os
import sys

sys.path.append("../functions")
from address_extractor import get_detailed_address

%matplotlib inline
```

Entrée [2]:

```python
get_detailed_address("Dhaka, Dhaka, Rampura")
```

Out[2]: {'City': 'Dhaka', 'Area': 'Rampura', 'Address': ''}

Entrée [3]:

```
1   # """
2   #  The functions below are developed by one of the teammate from Omdena
3   # """
4
5   # def get_detailed_address(address):
6   #     address = address.title()
7   #     address_dict = {"City": "", "Area": "", "Address": ""}
8   #     splitted_address = address.split(',')
9
10  #     for i in reversed(splitted_address):
11  #         if get_city_name(i.strip()):
12  #             address_dict["City"] = i.strip()
13  #             splitted_address.remove(i)
14  #         elif get_area_name(i.strip()):
15  #             address_dict["Area"] = i.strip()
16  #             splitted_address.remove(i)
17
18  #     address_dict["Address"] = ','.join(splitted_address)
19
20  #     return address_dict
21
22
23  # def get_city_name(name):
24  #     cities = ['Dhaka', 'Chattogram', 'Narayanganj City', 'Gazipur', 'Sylhet']
25
26  #     try:
27  #         cities.index(name)
28  #         return True
29
30  #     except:
31  #         return False
32
33
34  # def get_area_name(name):
35
36  #     areas = ['10 No. North Kattali Ward', '11 No. South Kattali Ward', '15 No. Bagmoniram Ward',
37  #              '16 No. Chawk Bazaar Ward', '22 No. Enayet Bazaar Ward', '29 No. West Madarbari Ward',
38  #              '30 No. East Madarbari Ward', '31 No. Alkoron Ward', '32 No. Andarkilla Ward',
39  #              '33 No. Firingee Bazaar Ward', '36 Goshail Danga Ward', '4 No Chandgaon Ward',
40  #              '7 No. West Sholoshohor Ward', '9 No. North Pahartali Ward', 'Adabor', 'Aftab Nagar', 'Agargaon',
41  #              'Ambarkhana', 'Badda', 'Bakalia', 'Banani', 'Banani Dohs', 'Banasree', 'Banglamotors', 'Bangshal',
```

```
42  #           'Baridhara', 'Baridhara Dohs', 'Bashabo', 'Bashundhara R-A', 'Bayazid', 'Cantonment', 'Chandra',
43  #           'Dakshin Khan', 'Demra', 'Dhanmondi', 'Double Mooring', 'Dumni', 'East Nasirabad', 'Eskaton', 'Fat
44  #           'Firojshah Colony', 'Gazipur Sadar Upazila', 'Gulistan', 'Gulshan', 'Halishahar', 'Hathazari', 'Ha
45  #           'Hazaribag', 'Ibrahimpur', 'Jalalabad Housing Society', 'Jamal Khan', 'Jatra Bari', 'Joar Sahara',
46  #           'Kachukhet', 'Kafrul', 'Kakrail', 'Kalabagan', 'Kalachandpur', 'Kamrangirchar', 'Kathalbagan',
47  #           'Kazir Dewri', 'Keraniganj', 'Khilgaon', 'Khilkhet', 'Khulshi', 'Kotwali', 'Kuril', 'Lal Khan Baza
48  #           'Lalbagh', 'Lalmatia', 'Maghbazar', 'Malibagh', 'Maniknagar', 'Mirpur', 'Mohakhali', 'Mohakhali Do
49  #           'Mohammadpur', 'Motijheel', 'Mugdapara', 'Muradpur', 'Nadda', 'Narayanganj', 'New Market', 'Niketa
50  #           'Nikunja', 'North Shahjahanpur', 'Panchlaish', 'Paribagh', 'Patenga', 'Purbachal', 'Railway Colony
51  #           'Rampura', 'Riaj Uddin Bazar', 'Sagorika Bscic Industrial Area', 'Savar', 'Shahbagh', 'Shahjahanpu
52  #           'Shantinagar', 'Shegunbagicha', 'Shiddheswari', 'Shiddhirganj', 'Sholokbahar', 'Shyamoli', 'Shyamp
53  #           'Sreepur', 'Sutrapur', 'Taltola', 'Tejgaon', 'Turag', 'Uttar Khan', 'Uttar Lalkhan', 'Uttara', 'Za
54  #           'Zindabazar']
55
56
57  #     try:
58  #         areas.index(name)
59  #         return True
60
61  #     except:
62  #         return False
```

Entrée [ ]: | 1

Entrée [ ]: | 1

Entrée [4]:
```
1  # CSV folders
2
3  raw_data_folder="../../../data/Raw_Data"
4  cleaned_data_folder="../../../data/CLeaned_Data"
5
6  bproperty_folder= f"{raw_data_folder}/bproperty_spider"
7  cleaned_bproperty_folder= f"{cleaned_data_folder}/bproperty"
```

Entrée [ ]: | 1

Entrée [ ]: 1

Entrée [5]:

```python
target_df_dic = {
    "area":[],  # value in float, in sqft; 1 Katha = 720 sqft (Thanks @Kausthab Dutta Phukan)
    "building_type":[],
    "building_nature": [], # originally named commercial_type; value will be either Commercial or Residential

    # splitted from location column
    "city": [],
    "locality": [],
    "address":[],
    #"country": [],
    #"municipality":[],
    #"district":[],
    #...
    #"otherZoneArea":[], # create new column for any new zone information, and keep collaborators informed


    "num_bath_rooms":[], # for Commercial properties, give 0 as value (since that make sense), not NaN
    "num_bed_rooms":[], # for Commercial properties, give 0 as value (since that make sense), not NaN

    # convert currencies to BDT : 1 lakh=100000 BDT, 1 crore=10000000 BDT, 1 Arab= 1000000000 BDT (Thanks @Al Mom
    "price": [],

    "property_description":[],
    "property_overview":[],

    "purpose":[], # Either Rent/Sale

    # retrieved from amenities column: assuming in sample 1 amenities has {"k1":"v1","k2":"v2"}
    #   and in sample 2 amenities has {"k3":"v3"}, we create new columns in the dataframe based on the keys of
    #   the dictionnaries
    "k1":[],
    "k2":[],
    "k3":[],

    # when any relevant column from other csv files is added, inform collaborators so that they follow the same p
}

target_df = pd.DataFrame(target_df_dic)
```

```
39  target_df.T
```

Out[5]:

| area |
|---|
| **building_type** |
| **building_nature** |
| **city** |
| **locality** |
| **address** |
| **num_bath_rooms** |
| **num_bed_rooms** |
| **price** |
| **property_description** |
| **property_overview** |
| **purpose** |
| **k1** |
| **k2** |
| **k3** |

Entrée [ ]:
```
1
```

Entrée [ ]:
```
1
```

# Assessing bproperty_spider_2023-04-09T19-44-07

Entrée [6]:
```
1  bproperty_df=pd.read_csv(f"{bproperty_folder}/bproperty_spider_2023-04-18T01-34-24.csv")
2  bproperty_df.head().T
```

Out[6]:

| | 0 | 1 | |
|---|---|---|---|
| **amenities** | {'Balcony or Terrace': 'yes', 'Flooring': 'yes... | {'View': 'yes', 'Parking Spaces': ' 1', 'Balco... | {'View': 'yes', 'Balcony or Terrace': 'ye |
| **area** | 1,185 sqft | 2,464 sqft | 1,14( |
| **building_type** | Apartment | Apartment | Apart |
| **commercial_type** | False | False | |
| **image_url** | https://images-cdn.bproperty.com/thumbnails/15... | https://images-cdn.bproperty.com/thumbnails/15... | https://ima cdn.bproperty.com/thumbnails |
| **location** | Khilgaon, Dhaka | Dhanmondi, Dhaka | Block TA, Section 6, Mirpur, D |
| **num_bath_rooms** | NaN | 4 Baths | |
| **num_bed_rooms** | 3 Beds | 3 Beds | 3 |
| **price** | 61 Lakh | 2.89 Crore | 75 |
| **property_description** | Grab This 1185 Sq Ft Beautiful Flat Is Vacant ... | A Vibrant 2464 Sq Ft Residential Flat For Sale... | 1140 Sq Ft Nicely Planned Apartme Avai |
| **property_overview** | This flat consists of facilities you can think... | Ready to move in somewhere with everything nea... | A spacious 1140 Square Feet apartme N |
| **property_url** | https://www.bproperty.com/en/property/details-... | https://www.bproperty.com/en/property/details-... | https://www.bproperty.com/en/property/deta |
| **purpose** | For Sale | For Sale | For |

Entrée [7]:
```
1  bproperty_df.shape
```

Out[7]: (17329, 13)

Entrée [ ]:  `1`

Entrée [8]:  `1 bproperty_df.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 17329 entries, 0 to 17328
Data columns (total 13 columns):
 #   Column                Non-Null Count   Dtype
---  ------                --------------   -----
 0   amenities             16438 non-null   object
 1   area                  17329 non-null   object
 2   building_type         17329 non-null   object
 3   commercial_type       17329 non-null   bool
 4   image_url             17312 non-null   object
 5   location              17329 non-null   object
 6   num_bath_rooms        5691 non-null    object
 7   num_bed_rooms         12646 non-null   object
 8   price                 17329 non-null   object
 9   property_description  17329 non-null   object
 10  property_overview     17329 non-null   object
 11  property_url          17329 non-null   object
 12  purpose               17329 non-null   object
dtypes: bool(1), object(12)
memory usage: 1.6+ MB
```

- `area` column should be decimal, not string (quality issue)
- Replace column name `commercial_type` by `building_nature` (or any relevant name), and change its values to `residential` or `commercial` accordingly. (quality issue)
- `location` is has concatened information: city, district, sector, etc. Those informations should be splitted in their relevant columns (column `city`, column `district`, ...). (tidiness issue)
- `num_bath_rooms` and `num_bed_rooms` should be decimal, no string. (quality issue)

Entrée [ ]:  `1`

Entrée [9]:
```
1  bproperty_df["price"].unique()
```

Out[9]:  array(['61 Lakh', '2.89 Crore', '75 Lakh', ..., '24 Crore', '7.3 Crore',
           '92.1 Lakh'], dtype=object)

- `price` content is not uniform accross the dataset. Some are in `Lakh`, other in `Crore`, etc... The unit used for the price should be uniformized. A special attention should be paid to the fact that there are `price` without unit (a solution need to be found for them). (quality issue)
- `price` should be decimal, not string

Entrée [ ]:
```
1
```

Entrée [10]:
```
1  bproperty_df["property_description"][0]
```

Out[10]:  'Grab This 1185 Sq Ft Beautiful Flat Is Vacant For Sale In Khilgaon'

Entrée [11]:
```
1  bproperty_df["property_description"][15]
```

Out[11]:  'Find 1522 Sq Ft Flat Available For Sale In Pallabi'

Entrée [ ]:
```
1
```

Entrée [12]:
```
1  bproperty_df["property_overview"][0]
```

Out[12]:  'This flat consists of facilities you can think of for a proper living standards. The amazing floor plan of this fla
         t ensures the buyer decision even more rightful for what they are pursuing for. The washroom comes with all the up-t
         o-date fixtures. For your delish cooking essentials, you are getting a convenient kitchen having a nice kitchen coun
         ter. You would also have refreshing balconies to spend some good family times in your morning and evening hours with
         your family. There are plenty of places nearby so going for recreational activities outside is an easy option to pic
         k as well. Make yourself a happy buyer by calling us about this beautiful apartment right away!'

Entrée [13]:
```
1   bproperty_df["property_overview"][150]
```

Out[13]: 'This plot is designed with your desired home in mind. It will not only offer you a pleasant lifestyle but also will offer you complete privacy and security to ensure that your family and guests have a safe and enjoyable time. It also offers you a very good payment plan. This beautiful plot is now available with full building.  If you like the plot, call us right away to know more.'
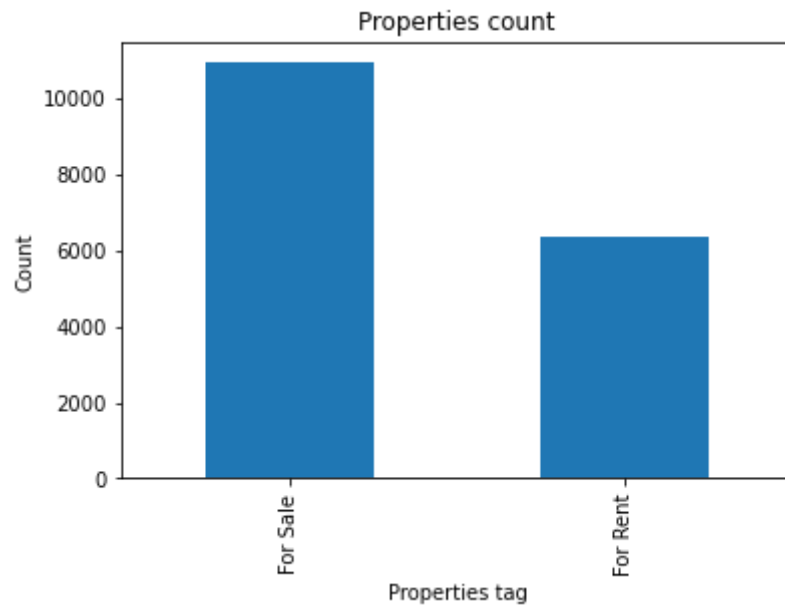
Entrée [ ]:
```
1
```

Entrée [14]:
```
1   property_per_purpose = bproperty_df["purpose"].value_counts()
2   property_per_purpose
```

Out[14]:
```
For Sale    10943
For Rent     6386
Name: purpose, dtype: int64
```

Entrée [15]:
```
1  property_per_purpose.plot(kind="bar")
2  plt.xlabel("Properties tag")
3  plt.ylabel("Count")
4  plt.title("Properties count");
```



Properties for sale are nearly the double of the properties for rent. And the amount of properties may be a little low to make the futures model predict well on unknow data.

Entrée [ ]:
```
1
```

Entrée [16]:
```
1  bproperty_df["amenities"][0]
```

Out[16]: "{'Balcony or Terrace': 'yes', 'Flooring': 'yes', 'Parking Spaces': ' 1', 'View': 'yes', 'Lobby in Building': 'yes', 'Electricity Backup': 'yes', 'Elevators in Building': ' 1', 'Floor Level': 'yes', 'CCTV Security': 'yes', 'Maintenance Staff': 'yes', 'Cleaning Services': 'yes'}"

Entrée [17]:
```
1  bproperty_df["amenities"][220]
```

Out[17]:  "{'Parking Spaces': ' 1', 'Flooring': 'yes', 'View': 'yes', 'Balcony or Terrace': 'yes', '24 Hours Concierge': 'yes', 'Double Glazed Windows': 'yes', 'Elevators in Building': ' 1', 'Floor Level': 'yes', 'Lobby in Building': 'yes', 'Broadband Internet': 'yes', 'Lawn or Garden': 'yes', 'CCTV Security': 'yes', 'Maintenance Staff': 'yes'}"

Each key in the dictionary of the feature `amenities` should become a column, with the following indications:

- `Floor level` : should be of type integer; its content should be the number of floor of the property
- `View` : should be of type boolean
- `Balcony or Terrace` : column should be named `balcony-or-terrace` , and should be of type boolean
- `Flooring` : should be of type boolean
- `Electricity backup` : column should be named `electricity-backup` , and should be of type boolean
- `Elevators in Buildings` : column should be named `elevator` , and should be of type int
- `Broadband Internet` : column should be named `internet` , and content should be boolean
- `CCTV Security` : column should be named `cctv-security` , and should be boolean
- `Cleaning Services` : column should be named `cleaning-services` , and should be boolean
- Keys present in the dictionary but not mentioned in the above list should also become a column

(tidiness issues)

Entrée [ ]:
```
1
```

Entrée [18]:
```
1  bproperty_df["building_type"].unique()
```

Out[18]:  array(['Apartment', 'Shop', 'Floor', 'Office', 'Building', 'Plot',
               'Duplex', 'Warehouse', 'Factory'], dtype=object)

Entrée [ ]:
```
1
```

Entrée [19]:
```
1  bproperty_df["purpose"].unique()
```

Out[19]:  array(['For Sale', 'For Rent'], dtype=object)

`purpose` should have `Rent` or `Sale` as values, to keep all cleaned datasets consistent.

Entrée [ ]:    1

**Assessment report summary**

**_Quality issues_**

1. `area` column should be decimal, not string.
2. Replace column name `commercial_type` by `building_nature`, and change its values to `residential` or `commercial` accordingly.
3. `num_bath_rooms` and `num_bed_rooms` should be decimal, no string.
4. `price` content is not uniform accross the dataset. Some are in `Lakh`, other in `Crore`, etc... The unit used for the price should be uniformized. Please pay attention to the fact that there are `price` without unit.
5. `price` should be decimal, not string
6. `purpose` should have `Rent` or `Sale` as values. This is not really an issue, its goal is only to keep values consistent accross all cleaned datasets.

**_Tidiness issues_**

1. `location` has concatened informations: city, district, sector, etc. Those informations will be splitted into `city` and `address`..
2. In `amenities` feature, each key in the dictionary should become a column, with the following indications:

   - `Floor level` : column should be named `floor-level`, and should be of type integer; its content should be the number of floor of the property ??
   - `View` : should be of type boolean
   - `Balcony or Terrace` : column should be named `balcony-or-terrace`, and should be of type boolean
   - `Flooring` : should be of type boolean
   - `Electricity backup` : column should be named `electricity-backup`, and should be of type boolean
   - `Elevators in Buildings` : column should be named `elevator`, and should be of type int
   - `Broadband Internet` : column should be named `internet`, and content should be boolean
   - `CCTV Security` : column should be named `cctv-security`, and should be boolean
   - `Cleaning Services` : column should be named `cleaning-services`, and should be boolean
   - Keys present in the dictionary but not mentioned in the above list should also become a column

Entrée [ ]:    1

Entrée [ ]:    1

# Cleaning bproperty

Entrée [ ]:    1

### `area` column should be decimal, not string ( [quality issue #1](quality%20issue%20%231) )

```
Entrée [20]:    1  # Recalling the type of area feature
                2  bproperty_df["area"].dtype
```

Out[20]:  dtype('O')

```
Entrée [21]:    1  bproperty_df["area"].unique()
```

Out[21]:  array(['1,185 sqft', '2,464 sqft', '1,140 sqft', ..., '3,842 sqft',
                  '13,350 sqft', '307 sqft'], dtype=object)

There are value in `sqft` and in `Katha`

**Define**

- Loop through `area` column, while:
  - converting `Katha` value to `sqft` value
  - removing the unit in the value, to only have the number left
- Convert `area` column to decimal

**Code**

Entrée [22]:

```python
"""
    Loop through `area` column, while:
        - converting `Katha` value to `sqft` value
        - removing the unit in the value, to only have the number left
"""

for index, row in bproperty_df.iterrows(): # loop through each sample

    # The code may take time, log in the console to keep track of things
    if index==0 or index%1000==0:
        print(f"Currently processing sample {index}...")

    # retrieve the area
    sample_area = bproperty_df.loc[index, "area"]
    splitted_sample_area = sample_area.split()

    # making sure there is only the value and the unit in sample_area
    if len(splitted_sample_area)>2:
        print(f"Sample of index {index} has a suspicious value as area: {sample_area}")
        break

    area = float( splitted_sample_area[0].replace(",","") ) # will contain the area; eg: 1345
    area_unit = splitted_sample_area[1].lower() # will contain the unit; eg: sqft

    # making sure all units are taken into account
    if area_unit not in ["sqft","katha"]:
        print(f"Sample of index {index} has a unit not taken into account for its area: {sample_area}")
        break

    # converting katha area to sqft area (1 Katha = 720 sqft => Thanks @Kausthab Dutta Phukan )
    if area_unit=="katha":
        area *= 720

    # updating the area of the sample in the dataframe
    bproperty_df.loc[index, "area"] = area

print("Processing has come to an end")

# Converting area to decimal
```

```
40  bproperty_df["area"] = bproperty_df["area"].astype(float)
```

```
Currently processing sample 0...
Currently processing sample 1000...
Currently processing sample 2000...
Currently processing sample 3000...
Currently processing sample 4000...
Currently processing sample 5000...
Currently processing sample 6000...
Currently processing sample 7000...
Currently processing sample 8000...
Currently processing sample 9000...
Currently processing sample 10000...
Currently processing sample 11000...
Currently processing sample 12000...
Currently processing sample 13000...
Currently processing sample 14000...
Currently processing sample 15000...
Currently processing sample 16000...
Currently processing sample 17000...
Processing has come to an end
```

Entrée [ ]:    1

### Testing

Entrée [23]:   1  bproperty_df["area"].dtype

Out[23]:  dtype('float64')

Entrée [ ]:    1

Entrée [ ]:    1

## Cleaning `commercial_type` feature ( [quality issue #2](quality issue #2) )

Replace column name `commercial_type` by `building_nature` , and change its values to `residential` or `commercial` accordingly.

Entrée [24]:
```
1  bproperty_df["commercial_type"].unique()
```

Out[24]: array([False,  True])

### Define

- Change column values: `True` is to be updated to `Commercial` , and `False` is to become `Residential`
- Replace column name ( `commercial_type` ) by `building_nature`

Entrée [ ]:
```
1
```

### Code

Entrée [25]:
```
1  # Replacing values of commercial_type column
2  bproperty_df.loc[ bproperty_df["commercial_type"]==True, ["commercial_type"] ] = "Commercial"
3  bproperty_df.loc[ bproperty_df["commercial_type"]==False, ["commercial_type"] ] = "Residential"
4
5  # Making sure values were updated
6  bproperty_df["commercial_type"].unique()
```

Out[25]: array(['Residential', 'Commercial'], dtype=object)

Entrée [26]:

```python
# Renaming column
bproperty_df.rename(columns={
    "commercial_type":"building_nature"
}, inplace=True)

# Confirming rename was done
bproperty_df.columns.to_list()
```

Out[26]: ['amenities',
 'area',
 'building_type',
 'building_nature',
 'image_url',
 'location',
 'num_bath_rooms',
 'num_bed_rooms',
 'price',
 'property_description',
 'property_overview',
 'property_url',
 'purpose']

Entrée [27]:
```
1  # Taking a look at content (for general confirmation)
2  bproperty_df.head(2).T
```

Out[27]:

| | 0 | 1 |
|---|---|---|
| **amenities** | {'Balcony or Terrace': 'yes', 'Flooring': 'yes... | {'View': 'yes', 'Parking Spaces': ' 1', 'Balco... |
| **area** | 1185.0 | 2464.0 |
| **building_type** | Apartment | Apartment |
| **building_nature** | Residential | Residential |
| **image_url** | https://images-cdn.bproperty.com/thumbnails/15... | https://images-cdn.bproperty.com/thumbnails/15... |
| **location** | Khilgaon, Dhaka | Dhanmondi, Dhaka |
| **num_bath_rooms** | NaN | 4 Baths |
| **num_bed_rooms** | 3 Beds | 3 Beds |
| **price** | 61 Lakh | 2.89 Crore |
| **property_description** | Grab This 1185 Sq Ft Beautiful Flat Is Vacant ... | A Vibrant 2464 Sq Ft Residential Flat For Sale... |
| **property_overview** | This flat consists of facilities you can think... | Ready to move in somewhere with everything nea... |
| **property_url** | https://www.bproperty.com/en/property/details-... | https://www.bproperty.com/en/property/details-... |
| **purpose** | For Sale | For Sale |

Entrée [ ]:
```
1
```

## `num_bath_rooms` and `num_bed_rooms` should be integer, no string. ( quality issue #3 )

Entrée [28]:
```
1  bproperty_df["num_bath_rooms"].dtype
```

Out[28]: dtype('O')

Entrée [29]:
```
1  bproperty_df["num_bath_rooms"].unique()
```

Out[29]: array([nan, '4 Baths', '3 Baths', '2 Baths', '10 Baths', '5 Baths',
          '8 Baths', '1 Bath', '7 Baths', '6 Baths', '9 Baths'], dtype=object)

Entrée [30]:
```
1  bproperty_df["num_bed_rooms"].dtype
```

Out[30]: dtype('O')

Entrée [31]:
```
1  bproperty_df["num_bed_rooms"].unique()
```

Out[31]: array(['3 Beds', '2 Beds', '4 Beds', nan, '21 Beds', '20 Beds', '5 Beds',
          '7 Beds', '1 Bed', '6 Beds', '19 Beds', '24 Beds', '33 Beds',
          '56 Beds', '18 Beds', '10 Beds', '13 Beds', '48 Beds', '12 Beds',
          '60 Beds', '40 Beds', '29 Beds', '23 Beds', '17 Beds', '14 Beds',
          '8 Beds', '50 Beds', '75 Beds', '42 Beds', '16 Beds', '36 Beds',
          '15 Beds', '25 Beds', '22 Beds', '46 Beds', '30 Beds', '11 Beds',
          '32 Beds', '94 Beds'], dtype=object)

Entrée [ ]:
```
1  
```

**Define**

- Replace `NaN` values by `0` (since in this case, that made sense: it mean the sample doesn't have a bath_room or bed_room
- Remove `Bed`, `Beds`, `Bath` and `Baths` from the values of `num_bed_rooms` and `num_bath_rooms`
- Convert `num_bed_rooms` and `num_bath_rooms` to integer

**Code**

Entrée [32]:
```
1  # Replace NaN value by 0 in num_bed_rooms and num_bath_rooms
2  bproperty_df["num_bed_rooms"].fillna("0", inplace=True)
3  bproperty_df["num_bath_rooms"].fillna("0", inplace=True)
4
5  # Check that NaN values where replaced
6  bproperty_df["num_bed_rooms"].isnull().sum(), bproperty_df["num_bath_rooms"].isnull().sum()
```

Out[32]:  (0, 0)

Entrée [33]:
```
1  # Removing the units (bed, bath, ...) in num_bed_rooms and num_bath_rooms
2  bproperty_df["num_bed_rooms"] = bproperty_df["num_bed_rooms"].apply(lambda x: x.split(" ")[0] )
3  bproperty_df["num_bath_rooms"] = bproperty_df["num_bath_rooms"].apply(lambda x: x.split(" ")[0] )
```

Entrée [34]:
```
1  # Converting num_bed_rooms and num_bath_rooms to integer
2  bproperty_df["num_bed_rooms"] = bproperty_df["num_bed_rooms"].astype(int)
3  bproperty_df["num_bath_rooms"] = bproperty_df["num_bath_rooms"].astype(int)
4
```

Entrée [ ]:
```
1
```

**Testing**

Entrée [35]:
```
1  # Checking type conversion was succesful
2  bproperty_df["num_bed_rooms"].dtype, bproperty_df["num_bath_rooms"].dtype
```

Out[35]:  (dtype('int32'), dtype('int32'))

Entrée [ ]:
```
1
```

Entrée [ ]:
```
1
```

## **`price` content is not uniform accross the dataset ( [quality issue #4 & #5](#) )**

`price` content is not uniform accross the dataset. Some are in `Lakh`, other in `Crore`, etc... The unit used for the price should be uniformized. A special attention should be paid to the fact that there are `price` without unit.

Furthermore, `price` should be decimal, not string.

Entrée [36]:
```
1  bproperty_df["price"].unique()
```

Out[36]:  array(['61 Lakh', '2.89 Crore', '75 Lakh', ..., '24 Crore', '7.3 Crore',
        '92.1 Lakh'], dtype=object)

**Define**

- Convert all price to the same currency
- Replace `Thousand` by triple `0`
- Convert the column to float

**Code**

Entrée [37]:

```python
1   """
2       Loop through `price` column, while:
3           * Converting all prices to BDT currency
4           * Replacing `Thousand` by triple `0`
5   """
6
7   for index, row in bproperty_df.iterrows(): # loop through each sample
8
9       # The code may take time, log in the console to keep track of things
10      if index==0 or index%1000==0:
11          print(f"Currently processing sample {index}...")
12
13      # retrieve the price
14      sample_price = bproperty_df.loc[index, "price"]
15      splitted_sample_price= sample_price.split()
16
17      # making sure there are only the value and unit in sample price
18      if len(splitted_sample_price)>2:
19          print(f"Sample of index {index} has a suspicious value as price: {sample_price}")
20          break
21
22      price = float( splitted_sample_price[0] ) # will contain the price; eg: 1345
23      price_unit = splitted_sample_price[1].lower() # will contain the unit; eg: Lakh, Crore
24
25      # making sure all units are taken into account
26      if price_unit not in ["arab","crore","lakh","thousand"]:
27          print(f"Sample of index {index} has a unit not taken into account for its price: {sample_price}")
28          break
29
30      # converting all price unit to BDT : 1 lakh=100000 BDT,1 crore=10000000 BDT, 1 Arab= 1000000000 BDT (Thanks @
31      if price_unit=="arab":
32          price *= 1000000000
33      elif price_unit=="crore":
34          price *= 10000000
35      elif price_unit=="lakh":
36          price *= 100000
37      elif price_unit=="thousand":
38          price *= 1000
39      else:
40          raise Exception(f"Currency {price_unit} not taken to account")
41
```

```python
42        # updating the price of the sample in the dataframe
43        bproperty_df.loc[index, "price"] = price
44
45 print("Processing has come to an end")
46
47 # Converting area to decimal
48 bproperty_df["price"] = bproperty_df["price"].astype(float)
```

```
Currently processing sample 0...
Currently processing sample 1000...
Currently processing sample 2000...
Currently processing sample 3000...
Currently processing sample 4000...
Currently processing sample 5000...
Currently processing sample 6000...
Currently processing sample 7000...
Currently processing sample 8000...
Currently processing sample 9000...
Currently processing sample 10000...
Currently processing sample 11000...
Currently processing sample 12000...
Currently processing sample 13000...
Currently processing sample 14000...
Currently processing sample 15000...
Currently processing sample 16000...
Currently processing sample 17000...
Processing has come to an end
```

Entrée [ ]:     1

**Testing**

Entrée [38]:    1 bproperty_df["price"].dtype

Out[38]: dtype('float64')

Entrée [ ]:  `1`

Entrée [ ]:  `1`

## Set `purpose` values to `Rent` or `Sale` ( [quality issue #6](#) )

`purpose` should have `Rent` or `Sale` as values. This is not really an issue, its goal is only to keep values consistent accross all cleaned datasets.

Entrée [39]:
```
1  bproperty_df["purpose"].unique()
```

Out[39]: `array(['For Sale', 'For Rent'], dtype=object)`

Entrée [ ]:  `1`

### Define

- Replace `For Sale` by `Sale`, and `For Rent` by `Rent`

Entrée [ ]:  `1`

### Code

Entrée [40]:
```
1  bproperty_df["purpose"] = bproperty_df["purpose"].apply(lambda x: x.split(" ")[1] )
```

### Testing

Entrée [41]:
```
1  bproperty_df["purpose"].unique()
```

Out[41]: `array(['Sale', 'Rent'], dtype=object)`

Entrée [ ]:  1

Entrée [ ]:  1

## Split `location` column content into adequate columns ( [tidiness issue #1](#) )

`location` has concatened informations: city, district, sector, etc. Those will be splitted into `city` and `address`.

Entrée [ ]:  1

Entrée [42]:  1  `bproperty_df["location"]`

Out[42]:
```
0                                    Khilgaon, Dhaka
1                                   Dhanmondi, Dhaka
2                   Block TA, Section 6, Mirpur, Dhaka
3                       Block J, Bashundhara R-A, Dhaka
4        Block M, South Banasree Project, Banasree, Dhaka
                             ...
17324                     Block H, Banasree, Dhaka
17325               Block J, Bashundhara R-A, Dhaka
17326               Block G, Bashundhara R-A, Dhaka
17327                         Baridhara DOHS, Dhaka
17328                      Block F, Banasree, Dhaka
Name: location, Length: 17329, dtype: object
```

Entrée [ ]:  1

**Define**

- Before
    - Split content of `location` to `city` and `address`
    - Remove `location` column
- Now
    - Retrieve the city, area, and address from each `location` through `get_detailed_address()`

- Update new columns (city, locality, address) based on values retrieve from `location`

Entrée [ ]:    ```
1
```

**Code**

Entrée [43]:
```
1  # testing the get_detailed_address
2  get_detailed_address(bproperty_df["location"][0])
```

Out[43]:  {'City': 'Dhaka', 'Area': 'Khilgaon', 'Address': ''}

Entrée [44]:
```
1  # bproperty_df["location"][13].title()
```

Entrée [45]:

```
1  # Create new columns
2  bproperty_df["city"] = np.NaN
3  bproperty_df["locality"] = np.NaN
4  bproperty_df["address"] = np.nan
5
6  # Check new columns
7  bproperty_df.head(3).T
```

Out[45]:

| | 0 | 1 |
|---|---|---|
| amenities | {'Balcony or Terrace': 'yes', 'Flooring': 'yes... | {'View': 'yes', 'Parking Spaces': ' 1', 'Balco... | {'View': 'yes', 'Balcony or Terrace': 'ye |
| area | 1185.0 | 2464.0 | 11 |
| building_type | Apartment | Apartment | Apart |
| building_nature | Residential | Residential | Reside |
| image_url | https://images-cdn.bproperty.com/thumbnails/15... | https://images-cdn.bproperty.com/thumbnails/15... | https://ima cdn.bproperty.com/thumbnails |
| location | Khilgaon, Dhaka | Dhanmondi, Dhaka | Block TA, Section 6, Mirpur, D |
| num_bath_rooms | 0 | 4 | |
| num_bed_rooms | 3 | 3 | |
| price | 6100000.0 | 28900000.0 | 75000 |
| property_description | Grab This 1185 Sq Ft Beautiful Flat Is Vacant ... | A Vibrant 2464 Sq Ft Residential Flat For Sale... | 1140 Sq Ft Nicely Planned Apartme Avai |
| property_overview | This flat consists of facilities you can think... | Ready to move in somewhere with everything nea... | A spacious 1140 Square Feet apartme N |
| property_url | https://www.bproperty.com/en/property/details-... | https://www.bproperty.com/en/property/details-... | https://www.bproperty.com/en/property/deta |
| purpose | Sale | Sale | |
| city | NaN | NaN | |
| locality | NaN | NaN | |
| address | NaN | NaN | |

Entrée [ ]: 
```
1
```

Entrée [46]:
```python
1  # Old code
2
3  # # Retrieve city in location
4  # bproperty_df["city"] = bproperty_df["location"].apply(lambda x: x.split(",")[-1].strip() )
5
6  # # Retrieve address in location
7  # bproperty_df["address"] = bproperty_df["location"].apply(lambda x: ",".join(x.split(",")[:-1]).strip() )
8
9  # # Checking the content of location, city, and address
10 # bproperty_df[ ["location","city","address"] ]
```

Entrée [47]:

```python
# New code

"""
    Loop through `location` column, while splitting each location to city, zone, address and add them
        to the relevant column
"""

for index, row in bproperty_df.iterrows(): # loop through each sample

    # The code may take time, log in the console to keep track of things
    if index==0 or index%1000==0:
        print(f"Currently processing sample {index}...")

    # retrieve the location
    location = bproperty_df.loc[index, "location"]

    # split the location to dictionary with Area, City, Address as keys
    location_dict = get_detailed_address(location)

    city = location_dict.get("City", np.NaN)
    locality = location_dict.get("Area", np.NaN)
    address = location_dict.get("Address", np.NaN)


    # updating the relevant columns of the sample in the dataframe
    bproperty_df.loc[index, "city"] = city
    bproperty_df.loc[index, "locality"] = locality
    bproperty_df.loc[index, "address"] = address

print("Processing has come to an end")
```

```
Currently processing sample 0...
Currently processing sample 1000...
Currently processing sample 2000...
Currently processing sample 3000...
Currently processing sample 4000...
Currently processing sample 5000...
Currently processing sample 6000...
Currently processing sample 7000...
Currently processing sample 8000...
Currently processing sample 9000...
Currently processing sample 10000...
Currently processing sample 11000...
Currently processing sample 12000...
Currently processing sample 13000...
Currently processing sample 14000...
Currently processing sample 15000...
Currently processing sample 16000...
Currently processing sample 17000...
Processing has come to an end
```

Entrée [48]:
```python
1  # Making sure the columns were splitted efficiently
2  bproperty_df[ ["location","city","locality","address"] ]
```

Out[48]:

| | location | city | locality | address |
|---|---|---|---|---|
| 0 | Khilgaon, Dhaka | Dhaka | Khilgaon | |
| 1 | Dhanmondi, Dhaka | Dhaka | Dhanmondi | |
| 2 | Block TA, Section 6, Mirpur, Dhaka | Dhaka | Mirpur | Block Ta, Section 6 |
| 3 | Block J, Bashundhara R-A, Dhaka | Dhaka | Bashundhara R-A | Block J |
| 4 | Block M, South Banasree Project, Banasree, Dhaka | Dhaka | Banasree | Block M, South Banasree Project |
| ... | ... | ... | ... | ... |
| 17324 | Block H, Banasree, Dhaka | Dhaka | Banasree | Block H |
| 17325 | Block J, Bashundhara R-A, Dhaka | Dhaka | Bashundhara R-A | Block J |
| 17326 | Block G, Bashundhara R-A, Dhaka | Dhaka | Bashundhara R-A | Block G |
| 17327 | Baridhara DOHS, Dhaka | Dhaka | Baridhara Dohs | |
| 17328 | Block F, Banasree, Dhaka | Dhaka | Banasree | Block F |

17329 rows × 4 columns

Entrée [49]:
```python
1  bproperty_df.shape
```

Out[49]:  (17329, 16)

Entrée [50]:
```python
1  # Drop location column
2  bproperty_df.drop(["location"], axis=1, inplace=True)
```

Entrée [51]:
```python
1  # Making sure removal was successful
2  bproperty_df.shape
```

Out[51]:  (17329, 15)

Entrée [ ]:    1

## Cleaning `amenities` feature ( [tidiness issue #2](#) )

In `amenities` feature, each key in the dictionaries (in its content) should become a column. The value of the key should become the sample value corresponding to that column.

Entrée [52]:    1  bproperty_df["amenities"][0]

Out[52]: "{'Balcony or Terrace': 'yes', 'Flooring': 'yes', 'Parking Spaces': ' 1', 'View': 'yes', 'Lobby in Building': 'yes', 'Electricity Backup': 'yes', 'Elevators in Building': ' 1', 'Floor Level': 'yes', 'CCTV Security': 'yes', 'Maintenance Staff': 'yes', 'Cleaning Services': 'yes'}"

Entrée [53]:    1  bproperty_df["amenities"][12]

Out[53]: "{'View': 'yes', 'Parking Spaces': ' 1', 'Floor Level': 'yes', 'Balcony or Terrace': 'yes', 'Lobby in Building': 'yes', 'Electricity Backup': 'yes', 'Flooring': 'yes', 'Elevators in Building': ' 1', 'Maintenance Staff': 'yes', 'Cleaning Services': 'yes'}"

Entrée [ ]:    1

**Define**

- Keys in the dictionaries of `amenities` will become new columns in the dataset; the values of the keys will become the new columns values for the corresponding sample.

Entrée [ ]:    1

**Code**

Entrée [54]:

```python
"""
    Loop through `amenities` column, while:
        * Converting the dictionnaries keys to new columns; the values of the keys are becoming
            the new columns values for the corresponding sample
"""

for index, row in bproperty_df.iterrows(): # loop through each sample

    # The code may take time, log in the console to keep track of things
    if index==0 or index%1000==0:
        print(f"Currently processing sample {index}...")

    # If current sample doesn't have amenities, go to the next one
    if pd.isna(bproperty_df.loc[index, "amenities"]):
        continue

    # retrieve the amenities
    sample_amenities = str(bproperty_df.loc[index, "amenities"]).replace("'","\"")

    amenities_dict = eval(sample_amenities)

    # Go through each key in the amenities dictionnary
    for key, value in amenities_dict.items():

        # put a suffix to the new column name, so that collaborators know it was generated from amenities feature
        column_name = slugify(key)+"-amenity"
        #print(column_name)

        # Create new column based on the key if not already existing
        if column_name not in bproperty_df.columns.to_list():
            bproperty_df[column_name]= np.NaN # Giving NaN as the default value for the column

        # Affecting to the new column created, for the current sample, the value of the dictionary's key
        bproperty_df.loc[index, column_name] = value
```

```
Currently processing sample 0...
Currently processing sample 1000...
Currently processing sample 2000...
Currently processing sample 3000...
Currently processing sample 4000...
Currently processing sample 5000...
Currently processing sample 6000...
Currently processing sample 7000...
Currently processing sample 8000...
Currently processing sample 9000...
Currently processing sample 10000...
Currently processing sample 11000...
Currently processing sample 12000...
Currently processing sample 13000...
Currently processing sample 14000...
Currently processing sample 15000...
Currently processing sample 16000...
Currently processing sample 17000...
```

Entrée [55]:

```
1  # Checking columns
2  bproperty_df.head(3).T
```

Out[55]:

| | 0 | 1 | |
|---|---|---|---|
| amenities | {'Balcony or Terrace': 'yes', 'Flooring': 'yes... | {'View': 'yes', 'Parking Spaces': ' 1', 'Balco... | {'View': 'yes', 'Balcony or Terrace': 'ye |
| area | 1185.0 | 2464.0 | 1 |
| building_type | Apartment | Apartment | Apart |
| building_nature | Residential | Residential | Resid |
| image_url | https://images-cdn.bproperty.com/thumbnails/15... | https://images-cdn.bproperty.com/thumbnails/15... | https://ima cdn.bproperty.com/thumbnails |
| num_bath_rooms | 0 | 4 | |
| num_bed_rooms | 3 | 3 | |
| price | 6100000.0 | 28900000.0 | 75000 |
| property_description | Grab This 1185 Sq Ft Beautiful Flat Is Vacant ... | A Vibrant 2464 Sq Ft Residential Flat For Sale... | 1140 Sq Ft Nicely Planned Apartme Avai |
| property_overview | This flat consists of facilities you can think... | Ready to move in somewhere with everything nea... | A spacious 1140 Square Feet apartme N |
| property_url | https://www.bproperty.com/en/property/details-... | https://www.bproperty.com/en/property/details-... | https://www.bproperty.com/en/property/deta |
| purpose | Sale | Sale | |
| city | Dhaka | Dhaka | D |
| locality | Khilgaon | Dhanmondi | N |
| address | | | Block Ta, Sect |
| balcony-or-terrace-amenity | yes | yes | |
| flooring-amenity | yes | yes | |
| parking-spaces-amenity | 1 | 1 | |
| view-amenity | yes | yes | |
| lobby-in-building-amenity | yes | yes | |
| electricity-backup-amenity | yes | yes | |

|  | 0 | 1 |
|---|---|---|
| elevators-in-building-amenity | 1 | 2 |
| floor-level-amenity | yes | yes |
| cctv-security-amenity | yes | yes |
| maintenance-staff-amenity | yes | NaN |
| cleaning-services-amenity | yes | NaN |
| service-elevators-amenity | NaN | yes |
| intercom-amenity | NaN | yes |
| atm-facility-amenity | NaN | yes |
| freehold-amenity | NaN | NaN |
| broadband-internet-amenity | NaN | NaN |
| double-glazed-windows-amenity | NaN | NaN |
| storage-areas-amenity | NaN | NaN |
| 24-hours-concierge-amenity | NaN | NaN |
| waste-disposal-amenity | NaN | NaN |
| lawn-or-garden-amenity | NaN | NaN |
| prayer-room-amenity | NaN | NaN |
| facilities-for-disabled-amenity | NaN | NaN |
| conference-room-amenity | NaN | NaN |
| furnished-amenity | NaN | NaN |

| | 0 | 1 |
|---|---|---|
| swimming-pool-amenity | NaN | NaN |
| steam-room-amenity | NaN | NaN |
| sauna-amenity | NaN | NaN |
| jacuzzi-amenity | NaN | NaN |
| barbeque-area-amenity | NaN | NaN |
| central-heating-amenity | NaN | NaN |
| business-center-amenity | NaN | NaN |
| first-aid-medical-center-amenity | NaN | NaN |
| day-care-center-amenity | NaN | NaN |
| shared-kitchen-amenity | NaN | NaN |
| cafeteria-or-canteen-amenity | NaN | NaN |
| laundry-facility-amenity | NaN | NaN |

Entrée [56]:
```
# Drop amenities column
bproperty_df.drop(["amenities"],axis=1, inplace=True)

# Check that removal was effective
"amenities" in bproperty_df.columns.to_list()
```

Out[56]: False

Entrée [ ]:     1

Entrée [ ]:     1

## Save cleaned dataset

Entrée [57]:
```python
1  # Create folder in which to save cleaned dataset
2  if not os.path.exists(cleaned_bproperty_folder):
3      os.makedirs(cleaned_bproperty_folder)
4      print(f"Create folder '{cleaned_bproperty_folder}'")
5  else:
6      print(f"Folder '{cleaned_bproperty_folder}' already exists")
```

Create folder '../../../data/CLeaned_Data/bproperty'

Entrée [58]:
```python
1  # Save cleaned dataset to csv
2  bproperty_df.to_csv(f"{cleaned_bproperty_folder}/cleaned_bproperty.csv", index=False)
```

Entrée [59]:

```
1  # Load saved csv (to make sure it was successfully save)
2  clean_bproperty_df = pd.read_csv(f"{cleaned_bproperty_folder}/cleaned_bproperty.csv")
3  clean_bproperty_df.head(3).T
```

```
C:\ProgramData\Anaconda3\lib\site-packages\IPython\core\interactiveshell.py:3165: DtypeWarning: Columns (16) have mi
xed types.Specify dtype option on import or set low_memory=False.
  has_raised = await self.run_ast_nodes(code_ast.body, cell_name,
```

localhost:8888/notebooks/Omdena/dhaka-bangladesh-real-estate-recommendation/src/tasks/task-2-data-preprocessing/bproperty -- cleaning/bproperty -- cleaning.ipynb#

46/49

Out[59]:

| | 0 | 1 | |
|---|---|---|---|
| area | 1185.0 | 2464.0 | 1 |
| building_type | Apartment | Apartment | Apart |
| building_nature | Residential | Residential | Reside |
| image_url | https://images-cdn.bproperty.com/thumbnails/15... | https://images-cdn.bproperty.com/thumbnails/15... | https://ima cdn.bproperty.com/thumbnails |
| num_bath_rooms | 0 | 4 | |
| num_bed_rooms | 3 | 3 | |
| price | 6100000.0 | 28900000.0 | 75000 |
| property_description | Grab This 1185 Sq Ft Beautiful Flat Is Vacant ... | A Vibrant 2464 Sq Ft Residential Flat For Sale... | 1140 Sq Ft Nicely Planned Apartme Avai |
| property_overview | This flat consists of facilities you can think... | Ready to move in somewhere with everything nea... | A spacious 1140 Square Feet apartme N |
| property_url | https://www.bproperty.com/en/property/details-... | https://www.bproperty.com/en/property/details-... | https://www.bproperty.com/en/property/deta |
| purpose | Sale | Sale | |
| city | Dhaka | Dhaka | D |
| locality | Khilgaon | Dhanmondi | N |
| address | NaN | NaN | Block Ta, Sect |
| balcony-or-terrace-amenity | yes | yes | |
| flooring-amenity | yes | yes | |
| parking-spaces-amenity | 1 | 1 | |
| view-amenity | yes | yes | |
| lobby-in-building-amenity | yes | yes | |
| electricity-backup-amenity | yes | yes | |
| elevators-in-building-amenity | 1.0 | 2.0 | |

localhost:8888/notebooks/Omdena/dhaka-bangladesh-real-estate-recommendation/src/tasks/task-2-data-preprocessing/bproperty -- cleaning/bproperty -- cleaning.ipynb#

47/49

| | 0 | 1 |
|---|---|---|
| floor-level-amenity | yes | yes |
| cctv-security-amenity | yes | yes |
| maintenance-staff-amenity | yes | NaN |
| cleaning-services-amenity | yes | NaN |
| service-elevators-amenity | NaN | yes |
| intercom-amenity | NaN | yes |
| atm-facility-amenity | NaN | yes |
| freehold-amenity | NaN | NaN |
| broadband-internet-amenity | NaN | NaN |
| double-glazed-windows-amenity | NaN | NaN |
| storage-areas-amenity | NaN | NaN |
| 24-hours-concierge-amenity | NaN | NaN |
| waste-disposal-amenity | NaN | NaN |
| lawn-or-garden-amenity | NaN | NaN |
| prayer-room-amenity | NaN | NaN |
| facilities-for-disabled-amenity | NaN | NaN |
| conference-room-amenity | NaN | NaN |
| furnished-amenity | NaN | NaN |
| swimming-pool-amenity | NaN | NaN |

| | 0 | 1 |
|---|---|---|
| steam-room-amenity | NaN | NaN |
| sauna-amenity | NaN | NaN |
| jacuzzi-amenity | NaN | NaN |
| barbeque-area-amenity | NaN | NaN |
| central-heating-amenity | NaN | NaN |
| business-center-amenity | NaN | NaN |
| first-aid-medical-center-amenity | NaN | NaN |
| day-care-center-amenity | NaN | NaN |
| shared-kitchen-amenity | NaN | NaN |
| cafeteria-or-canteen-amenity | NaN | NaN |
| laundry-facility-amenity | NaN | NaN |

Entrée [ ]:    1

Entrée [ ]:    1

Entrée [ ]:    1