



## Homework 03

**Name: Miran Tafazzul Hussain Junaidi**

**Student ID: 3034487132**

**Note: Please print the output of each question in a new cell below your code**

### Numpy Introduction

**1a) Create two numpy arrays (a and b). a should be all integers between 25-34 (inclusive), and b should be ten evenly spaced numbers between 1-6. Print all the results below and store them separately:**

- i) Cube (i.e. raise to the power of 3) all the elements in both arrays (element-wise)
- ii) Add both the cubed arrays (e.g.,  $[1,2] + [3,4] = [4,6]$ )
- iii) Sum the elements with even indices of the added array.
- iv) Take the square root of the added array (element-wise square root)

In [1]:

```
# your code here
#creating the array a
import numpy as np
a = np.array([x for x in range(25,35)])
a
```

Out[1]:

```
array([25, 26, 27, 28, 29, 30, 31, 32, 33, 34])
```

In [2]:

```
#Creating the array b
b_list= [1]
for i in range (9):
    b_list.append(b_list[i]+ 5/10)
b = np.array(b_list)
b
```

Out[2]:

```
array([1. , 1.5, 2. , 2.5, 3. , 3.5, 4. , 4.5, 5. , 5.5])
```

In [3]:

```
#Cubing both the arrays
a_cubed = a**3
b_cubed = b**3
print(a_cubed)
print(b_cubed)
```

```
[15625 17576 19683 21952 24389 27000 29791 32768 35937 39304]
[  1.    3.375  8.    15.625 27.    42.875 64.    91.125 125.
 166.375]
```

In [4]:

```
cubed_sum = np.add(a_cubed,b_cubed)
cubed_sum
```

Out[4]:

```
array([15626.    , 17579.375, 19691.    , 21967.625, 24416.    , 27042.87
5,
      29855.    , 32859.125, 36062.    , 39470.375])
```

In [5]:

```
#Sum of the elements with even indices of the cubed and added arrays
even_index_sum=0

for i in cubed_sum[::2] :
    even_index_sum += i
even_index_sum
```

Out[5]:

```
125650.0
```

In [6]:

```
# Square root of cubed_sum
cbd_sum_sqrt = cubed_sum **0.5
cbd_sum_sqrt
```

Out[6]:

```
array([125.00399994, 132.58723543, 140.32462364, 148.21479346,
      156.25619988, 164.44717997, 172.7859948 , 181.27086087,
      189.89997367, 198.67152539])
```

**1b) Append b to a, reshape the appended array so that it is a 4x5, 2d array and store the results in a variable called m. Print m.**

In [7]:

```
# your code here
m = np.concatenate([a,b], axis = 0)
m = m.reshape(4,5)
m
```

Out[7]:

```
array([[25. , 26. , 27. , 28. , 29. ],
       [30. , 31. , 32. , 33. , 34. ],
       [ 1. ,  1.5,  2. ,  2.5,  3. ],
       [ 3.5,  4. ,  4.5,  5. ,  5.5]])
```

**1c) Extract the third and the fourth column of the m matrix. Store the resulting 4x2 matrix in a new variable called m2. Print m2.**

In [8]:

```
# your code here
m2 = m[:,2:4]
m2
```

Out[8]:

```
array([[27. , 28. ],
       [32. , 33. ],
       [ 2. ,  2.5],
       [ 4.5,  5. ]])
```

**1d) Take the dot product of m2 and m store the results in a matrix called m3. Print m3. Note that Dot product of two matrices  $A.B = A^T B$**

In [9]:

```
# your code here
m2_t = m2.transpose()
m3 = m2_t.dot(m)
m3
```

Out[9]:

```
array([[1652.75, 1715. , 1777.25, 1839.5 , 1901.75],
       [1710. , 1774.75, 1839.5 , 1904.25, 1969. ]])
```

## Numpy conditions

**2a) Create a numpy array called 'f' where the values are cosine(x) for x from 0 to pi with 50 equally spaced values in f**

- Print f
- Use condition on the array and print an array that is True when  $f \geq 1/2$  and False when  $f < 1/2$
- Create and print an array sequence that has only those values where  $f \geq 1/2$

In [10]:

```
# your code here
import math
f_noncos = np.linspace(0,np.pi,50)
f = np.cos(f_noncos)
f
```

Out[10]:

```
array([ 1.          ,  0.99794539,  0.99179001,  0.98155916,  0.9672948
6,
        0.94905575,  0.92691676,  0.90096887,  0.8713187 ,  0.8380881
,
        0.80141362,  0.76144596,  0.71834935,  0.67230089,  0.6234898
,
        0.57211666,  0.51839257,  0.46253829,  0.40478334,  0.3453650
5,
        0.28452759,  0.22252093,  0.1595999 ,  0.09602303,  0.0320515
8,
       -0.03205158, -0.09602303, -0.1595999 , -0.22252093, -0.2845275
9,
       -0.34536505, -0.40478334, -0.46253829, -0.51839257, -0.5721166
6,
       -0.6234898 , -0.67230089, -0.71834935, -0.76144596, -0.8014136
2,
       -0.8380881 , -0.8713187 , -0.90096887, -0.92691676, -0.9490557
5,
       -0.96729486, -0.98155916, -0.99179001, -0.99794539, -1.
])
```

In [12]:

```
#Boolean array
f_bool = f>=0.5
f_bool
```

Out[12]:

```
array([ True,  True,  True,  True,  True,  True,  True,  True,  True,
        True,  True,  True,  True,  True,  True,  True,  True, False,
        False, False, False, False, False, False, False, False, False,
        False, False, False, False, False, False, False, False, False,
        False, False, False, False, False, False, False, False, False,
        False, False, False, False, False])
```

In [13]:

```
#f with elements grater than or equal to 1/2
f[f_bool]
```

Out[13]:

```
array([1.          ,  0.99794539,  0.99179001,  0.98155916,  0.96729486,
        0.94905575,  0.92691676,  0.90096887,  0.8713187 ,  0.8380881 ,
        0.80141362,  0.76144596,  0.71834935,  0.67230089,  0.6234898 ,
        0.57211666,  0.51839257])
```

## NumPy and 2 Variable Prediction

Let 'x' be the number of miles a person drives per day and 'y' be the dollars spent on buying car fuel (per day).

We have created 2 numpy arrays each of size 100 that represent x and y.

x ( number of miles) ranges from 1 to 10 with a uniform noise of (0,1/2)

y (money spent in dollars) will be from 1 to 20 with a uniform noise (0,1)

In [14]:

```
# seed the random number generator with a fixed value
import numpy as np
np.random.seed(500)

x = np.linspace(1,10,100)+ np.random.uniform(low=0,high=.5,size=100)
y = np.linspace(1,20,100)+ np.random.uniform(low=0,high=1,size=100)
print('x = ',x)
print('y = ',y)
```

x = [ 1.34683976 1.12176759 1.51512398 1.55233174 1.40619168 1.6  
5075498  
1.79399331 1.80243817 1.89844195 2.00100023 2.3344038 2.224248  
72  
2.24914511 2.36268477 2.49808849 2.8212704 2.68452475 2.682294  
27  
3.09511169 2.95703884 3.09047742 3.2544361 3.41541904 3.408863  
75  
3.50672677 3.74960644 3.64861355 3.7721462 3.56368566 4.010927  
01  
4.15630694 4.06088549 4.02517179 4.25169402 4.15897504 4.268353  
33  
4.32520644 4.48563164 4.78490721 4.84614839 4.96698768 5.187542  
59  
5.29582013 5.32097781 5.0674106 5.47601124 5.46852704 5.645374  
52  
5.49642807 5.89755027 5.68548923 5.76276141 5.94613234 6.181357  
13  
5.96522091 6.0275473 6.54290191 6.4991329 6.74003765 6.818098  
07  
6.50611821 6.91538752 7.01250925 6.89905417 7.31314433 7.204722  
97  
7.1043621 7.48199528 7.58957227 7.61744354 7.6991707 7.854368  
22  
8.03510784 7.80787781 8.22410224 7.99366248 8.40581097 8.289137  
92  
8.45971515 8.54227144 8.6906456 8.61856507 8.83489887 8.663096  
58  
8.94837987 9.20890222 8.9614749 8.92608294 9.13231416 9.558898  
96  
9.61488451 9.54252979 9.42015491 9.90952569 10.00659591 10.025042  
65  
10.07330937 9.93489915 10.0892334 10.36509991]  
y = [ 1.6635012 2.0214592 2.10816052 2.26016496 1.96287558 2.9  
554635  
3.02881887 3.33565296 2.75465779 3.4250107 3.39670148 3.393777  
67  
3.78503343 4.38293049 4.32963586 4.03925039 4.73691868 4.300983  
99  
4.8416329 4.78175957 4.99765787 5.31746817 5.76844671 5.937237  
49  
5.72811642 6.70973615 6.68143367 6.57482731 7.17737603 7.548632  
52  
7.30221419 7.3202573 7.78023884 7.91133365 8.2765417 8.692032  
81  
8.78219865 8.45897546 8.89094715 8.81719921 8.87106971 9.661925  
62  
9.4020625 9.85990783 9.60359778 10.07386266 10.6957995 10.667219  
16  
11.18256285 10.57431836 11.46744716 10.94398916 11.26445259 12.097548

```

28
12.11988037 12.121557 12.17613693 12.43750193 13.00912372 12.864071
94
13.24640866 12.76120085 13.11723062 14.07841099 14.19821707 14.272890
01
14.30624942 14.63060835 14.2770918 15.0744923 14.45261619 15.118973
13
15.2378667 15.27203124 15.32491892 16.01095271 15.71250558 16.294885
06
16.70618934 16.56555394 16.42379457 17.18144744 17.13813976 17.696136
25
17.37763019 17.90942839 17.90343733 18.01951169 18.35727914 18.168412
69
18.61813748 18.66062754 18.81217983 19.44995194 19.7213867 19.719667
26
19.78961904 19.64385088 20.69719809 20.07974319]

```

**\_\_3a) Find Expected value of x and the expected value of y\_\_**

In [15]:

```

# your code here
e_x = np.mean(x)
e_y = np.mean(y)
print("Expected Value of X = " , e_x)
print("Expected Value of Y = " , e_y)

```

```

Expected Value of X = 5.782532541587923
Expected Value of Y = 11.012981683344968

```

**3b) Find variance of distributions of x and y**

In [16]:

```

# your code here
print("Variance of X = ", np.var(x))
print("Variance of Y = ", np.var(y))

```

```

Variance of X = 7.03332752947585
Variance of Y = 30.113903575509635

```

**3c) Find co-variance of x and y.**

In [17]:

```

# your code here
xy = x*y
e_xy = np.mean(xy)
covar = e_xy - e_x * e_y
covar

```

Out[17]:

```
14.511166394475424
```

**3d) Assuming that number of dollars spent in car fuel is only dependant on the miles driven, by a linear relationship.**

**Write code that uses a linear predictor to calculate a predicted value of y for each x i.e  $y_{\text{predicted}} = f(x)$**

**=  $y_0 + mx$ . (Do not use Machine learning libraries)**

In [18]:

```
# your code here
m = covar / np.var(x)
c = e_y - m*e_x
m
x_in = int(input("Enter the miles driven"))
def pred_cost(miles,m,c):
    y_ou = m*x_in + c
    return y_ou
print("The amount spent on feul is $", pred_cost(x_in, m, c))
```

Enter the miles driven10

The amount spent on feul is \$ 19.71446356312654

**3e) Predict y for each value in x, put the error into an array called y\_error**



In [19]:

```
# your code here
y_p = m*x + c
y_error = y_p - y
y_error
```

Out[19]:

```
array([ 0.19775597, -0.62457111,  0.10030076,  0.02506341,  0.0208364
9,
      -0.46716823, -0.24499418, -0.53440482,  0.24466541, -0.2140891
8,
      0.50209852,  0.27775029, -0.06213923, -0.42578118, -0.0931215
,
      0.86405311, -0.1157489 ,  0.31558388,  0.62666017,  0.4016614
9,
      0.46107377,  0.47954311,  0.3607047 ,  0.17838904,  0.5894211
6,
      0.10891094, -0.07115518,  0.29032384, -0.74232081, -0.1908286
3,
      0.35553767,  0.14062095, -0.39304511, -0.0567791 , -0.6132850
2,
     -0.80310676, -0.77597321, -0.12176065,  0.06373323,  0.2638340
2,
      0.45927925,  0.12347238,  0.60673379,  0.20079382, -0.0660562
,
      0.30670405, -0.33067419,  0.062778 , -0.75987212,  0.6759679
8,
     -0.65468531,  0.02820071,  0.08606832, -0.26171143, -0.7299759
2,
     -0.60306068,  0.40563939,  0.05397013, -0.02061681,  0.2854892
8,
     -0.7405245 ,  0.58908804,  0.43343988, -0.76182107, -0.0272760
4,
     -0.32564401, -0.56606805, -0.11129392,  0.46417555, -0.2757209
3,
      0.5147747 ,  0.16862142,  0.42262995, -0.08035574,  0.7255111
2,
     -0.43596616,  0.71282602, -0.11027337, -0.16964259,  0.1413230
1,
      0.58920807, -0.31716141,  0.17248631, -0.73997278,  0.1671299
7,
      0.17284167, -0.33165948, -0.52075457, -0.43302563,  0.6359709
,
      0.30175553,  0.10998314, -0.29405306,  0.07784496,  0.0066855
4,
      0.04646431,  0.07609646, -0.06370343, -0.79862812,  0.3879947
7])
```

**3f) Write code that calculates the root mean square error(RMSE), that is root of average of y-error squared**

In [20]:

```
#your code here
RMS_Error = np.sqrt((y_error**2).mean())
RMS_Error
```

Out[20]:

0.41767772366856115

## Pandas Introduction

### Reading File

In [21]:

```
# Load required modules
import pandas as pd
import numpy as np
```

Read the CSV file called 'data3.csv' into a dataframe called df.

#### Data description

- File location: [https://bcourses.berkeley.edu/files/74463396/download?download\\_frd=1](https://bcourses.berkeley.edu/files/74463396/download?download_frd=1)  
([https://bcourses.berkeley.edu/files/74463396/download?download\\_frd=1](https://bcourses.berkeley.edu/files/74463396/download?download_frd=1))
- Data source: [http://www.fao.org/nr/water/aquastat/data/query/index.html?\\*lang=en](http://www.fao.org/nr/water/aquastat/data/query/index.html?*lang=en)  
([http://www.fao.org/nr/water/aquastat/data/query/index.html?\\*lang=en](http://www.fao.org/nr/water/aquastat/data/query/index.html?*lang=en))
- Data, units:
- GDP, current USD (CPI adjusted)
- NRI, mm/yr
- Population density, inhab/km<sup>2</sup>
- Total area of the country, 1000 ha = 10km<sup>2</sup>
- Total Population, unit 1000 inhabitants

In [22]:

```
# your code here
data = pd.read_csv("data3.csv")
data
```

Out[22]:

	Area	Area Id	Variable Name	Variable Id	Year	Value	Symbol	Other
0	Argentina	9.0	Total area of the country	4100.0	1962.0	2.780400e+05	E	NaN
1	Argentina	9.0	Total area of the country	4100.0	1967.0	2.780400e+05	E	NaN
2	Argentina	9.0	Total area of the country	4100.0	1972.0	2.780400e+05	E	NaN
3	Argentina	9.0	Total area of the country	4100.0	1977.0	2.780400e+05	E	NaN
4	Argentina	9.0	Total area of the country	4100.0	1982.0	2.780400e+05	E	NaN
5	Argentina	9.0	Total area of the country	4100.0	1987.0	2.780400e+05	E	NaN
6	Argentina	9.0	Total area of the country	4100.0	1992.0	2.780400e+05	E	NaN
7	Argentina	9.0	Total area of the country	4100.0	1997.0	2.780400e+05	E	NaN
8	Argentina	9.0	Total area of the country	4100.0	2002.0	2.780400e+05	E	NaN
9	Argentina	9.0	Total area of the country	4100.0	2007.0	2.780400e+05	E	NaN
10	Argentina	9.0	Total area of the country	4100.0	2012.0	2.780400e+05	E	NaN
11	Argentina	9.0	Total area of the country	4100.0	2014.0	2.780400e+05	E	NaN
12	Argentina	9.0	Total population	4104.0	1962.0	2.128800e+04	E	NaN
13	Argentina	9.0	Total population	4104.0	1967.0	2.293200e+04	E	NaN
14	Argentina	9.0	Total population	4104.0	1972.0	2.478300e+04	E	NaN
15	Argentina	9.0	Total population	4104.0	1977.0	2.687900e+04	E	NaN
16	Argentina	9.0	Total population	4104.0	1982.0	2.899400e+04	E	NaN
17	Argentina	9.0	Total population	4104.0	1987.0	3.132600e+04	E	NaN
18	Argentina	9.0	Total population	4104.0	1992.0	3.365500e+04	E	NaN
19	Argentina	9.0	Total population	4104.0	1997.0	3.583400e+04	E	NaN
20	Argentina	9.0	Total population	4104.0	2002.0	3.788900e+04	E	NaN

	Area	Area Id	Variable Name	Variable Id	Year	Value	Symbol	Other
21	Argentina	9.0	Total population	4104.0	2007.0	3.997000e+04	E	NaN
22	Argentina	9.0	Total population	4104.0	2012.0	4.209500e+04	E	NaN
23	Argentina	9.0	Total population	4104.0	2015.0	4.341700e+04	E	NaN
24	Argentina	9.0	Population density	4107.0	1962.0	7.656000e+00	E	NaN
25	Argentina	9.0	Population density	4107.0	1967.0	8.248000e+00	E	NaN
26	Argentina	9.0	Population density	4107.0	1972.0	8.913000e+00	E	NaN
27	Argentina	9.0	Population density	4107.0	1977.0	9.667000e+00	E	NaN
28	Argentina	9.0	Population density	4107.0	1982.0	1.043000e+01	E	NaN
29	Argentina	9.0	Population density	4107.0	1987.0	1.127000e+01	E	NaN
...	...	...	...	...	...	...	...	...
368	United States of America	231.0	Population density	4107.0	2012.0	3.202000e+01	E	NaN
369	United States of America	231.0	Population density	4107.0	2015.0	3.273000e+01	E	NaN
370	United States of America	231.0	Gross Domestic Product (GDP)	4112.0	1962.0	6.050000e+11	E	NaN
371	United States of America	231.0	Gross Domestic Product (GDP)	4112.0	1967.0	8.620000e+11	E	NaN
372	United States of America	231.0	Gross Domestic Product (GDP)	4112.0	1972.0	1.280000e+12	E	NaN
373	United States of America	231.0	Gross Domestic Product (GDP)	4112.0	1977.0	2.090000e+12	E	NaN
374	United States of America	231.0	Gross Domestic Product (GDP)	4112.0	1982.0	3.340000e+12	E	NaN
375	United States of America	231.0	Gross Domestic Product (GDP)	4112.0	1987.0	4.870000e+12	E	NaN
376	United States of America	231.0	Gross Domestic Product (GDP)	4112.0	1992.0	6.540000e+12	E	NaN
377	United States of America	231.0	Gross Domestic Product (GDP)	4112.0	1997.0	8.610000e+12	E	NaN
378	United States of America	231.0	Gross Domestic Product (GDP)	4112.0	2002.0	1.100000e+13	E	NaN

	Area	Area Id	Variable Name	Variable Id	Year	Value	Symbol	Other
379	United States of America	231.0	Gross Domestic Product (GDP)	4112.0	2007.0	1.450000e+13	E	NaN
380	United States of America	231.0	Gross Domestic Product (GDP)	4112.0	2012.0	1.620000e+13	E	NaN
381	United States of America	231.0	Gross Domestic Product (GDP)	4112.0	2015.0	1.790000e+13	E	NaN
382	United States of America	231.0	National Rainfall Index (NRI)	4472.0	1965.0	9.285000e+02	E	NaN
383	United States of America	231.0	National Rainfall Index (NRI)	4472.0	1969.0	9.522000e+02	E	NaN
384	United States of America	231.0	National Rainfall Index (NRI)	4472.0	1974.0	1.008000e+03	E	NaN
385	United States of America	231.0	National Rainfall Index (NRI)	4472.0	1981.0	9.492000e+02	E	NaN
386	United States of America	231.0	National Rainfall Index (NRI)	4472.0	1984.0	9.746000e+02	E	NaN
387	United States of America	231.0	National Rainfall Index (NRI)	4472.0	1992.0	1.020000e+03	E	NaN
388	United States of America	231.0	National Rainfall Index (NRI)	4472.0	1996.0	1.005000e+03	E	NaN
389	United States of America	231.0	National Rainfall Index (NRI)	4472.0	2002.0	9.387000e+02	E	NaN
390	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
391	E - External data	NaN	NaN	NaN	NaN	NaN	NaN	NaN
392	I - AQUASTAT estimate	NaN	NaN	NaN	NaN	NaN	NaN	NaN
393	K - Aggregate data	NaN	NaN	NaN	NaN	NaN	NaN	NaN
394	L - Modelled data	NaN	NaN	NaN	NaN	NaN	NaN	NaN
395	(c) FAO of the UN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
396	The information contained in AQUASTAT is provi...	NaN	NaN	NaN	NaN	NaN	NaN	NaN
397	FAO. 2016. AQUASTAT Main Database - Food and A...	NaN	NaN	NaN	NaN	NaN	NaN	NaN

398 rows × 8 columns

**4a ) Display the first 10 rows of the dataframe**

In [23]:

```
# your code here
data[0:10]
```

Out[23]:

	Area	Area Id	Variable Name	Variable Id	Year	Value	Symbol	Other
0	Argentina	9.0	Total area of the country	4100.0	1962.0	278040.0	E	NaN
1	Argentina	9.0	Total area of the country	4100.0	1967.0	278040.0	E	NaN
2	Argentina	9.0	Total area of the country	4100.0	1972.0	278040.0	E	NaN
3	Argentina	9.0	Total area of the country	4100.0	1977.0	278040.0	E	NaN
4	Argentina	9.0	Total area of the country	4100.0	1982.0	278040.0	E	NaN
5	Argentina	9.0	Total area of the country	4100.0	1987.0	278040.0	E	NaN
6	Argentina	9.0	Total area of the country	4100.0	1992.0	278040.0	E	NaN
7	Argentina	9.0	Total area of the country	4100.0	1997.0	278040.0	E	NaN
8	Argentina	9.0	Total area of the country	4100.0	2002.0	278040.0	E	NaN
9	Argentina	9.0	Total area of the country	4100.0	2007.0	278040.0	E	NaN

**4b) Display the column names.**

In [24]:

```
# your code here
list(data)
```

Out[24]:

```
['Area',
 'Area Id',
 'Variable Name',
 'Variable Id',
 'Year',
 'Value',
 'Symbol',
 'Other']
```

**4c) Use iloc to display the first 3 rows and first 4 columns.**

In [25]:

```
# your code here
data.iloc[0:3,0:4]
```

Out[25]:

	Area	Area Id	Variable Name	Variable Id
0	Argentina	9.0	Total area of the country	4100.0
1	Argentina	9.0	Total area of the country	4100.0
2	Argentina	9.0	Total area of the country	4100.0

## Data Preprocessing

5a ) Find all the rows that have 'NaN' in the 'Symbol' column. Display first 5 rows.

*Hint : You might have to use a condition (mask)*

In [26]:

```
# your code here

No_sym_df = data[:,data['Symbol'] != 'E' ]
No_sym_df[0:5]
```

Out[26]:

	Area	Area Id	Variable Name	Variable Id	Year	Value	Symbol	Other
390		NaN	NaN	NaN	NaN	NaN	NaN	NaN
391	E - External data	NaN	NaN	NaN	NaN	NaN	NaN	NaN
392	I - AQUASTAT estimate	NaN	NaN	NaN	NaN	NaN	NaN	NaN
393	K - Aggregate data	NaN	NaN	NaN	NaN	NaN	NaN	NaN
394	L - Modelled data	NaN	NaN	NaN	NaN	NaN	NaN	NaN

5b ) Now, we will try to get rid of the NaN valued rows and columns. First, drop the column 'Other' which only has 'NaN' values. Then drop all other rows that have any column with a value 'NaN'. Then display the last 5 rows of the dataframe.

In [27]:

```
# your code here
data_no_nan = data.drop("Other", axis = 1)
data_no_nan = data_no_nan.dropna()
data_no_nan
```

Out[27]:

	Area	Area Id	Variable Name	Variable Id	Year	Value	Symbol
0	Argentina	9.0	Total area of the country	4100.0	1962.0	2.780400e+05	E
1	Argentina	9.0	Total area of the country	4100.0	1967.0	2.780400e+05	E
2	Argentina	9.0	Total area of the country	4100.0	1972.0	2.780400e+05	E
3	Argentina	9.0	Total area of the country	4100.0	1977.0	2.780400e+05	E
4	Argentina	9.0	Total area of the country	4100.0	1982.0	2.780400e+05	E
5	Argentina	9.0	Total area of the country	4100.0	1987.0	2.780400e+05	E
6	Argentina	9.0	Total area of the country	4100.0	1992.0	2.780400e+05	E
7	Argentina	9.0	Total area of the country	4100.0	1997.0	2.780400e+05	E
8	Argentina	9.0	Total area of the country	4100.0	2002.0	2.780400e+05	E
9	Argentina	9.0	Total area of the country	4100.0	2007.0	2.780400e+05	E
10	Argentina	9.0	Total area of the country	4100.0	2012.0	2.780400e+05	E
11	Argentina	9.0	Total area of the country	4100.0	2014.0	2.780400e+05	E
12	Argentina	9.0	Total population	4104.0	1962.0	2.128800e+04	E
13	Argentina	9.0	Total population	4104.0	1967.0	2.293200e+04	E
14	Argentina	9.0	Total population	4104.0	1972.0	2.478300e+04	E
15	Argentina	9.0	Total population	4104.0	1977.0	2.687900e+04	E
16	Argentina	9.0	Total population	4104.0	1982.0	2.899400e+04	E
17	Argentina	9.0	Total population	4104.0	1987.0	3.132600e+04	E
18	Argentina	9.0	Total population	4104.0	1992.0	3.365500e+04	E
19	Argentina	9.0	Total population	4104.0	1997.0	3.583400e+04	E
20	Argentina	9.0	Total population	4104.0	2002.0	3.788900e+04	E
21	Argentina	9.0	Total population	4104.0	2007.0	3.997000e+04	E
22	Argentina	9.0	Total population	4104.0	2012.0	4.209500e+04	E
23	Argentina	9.0	Total population	4104.0	2015.0	4.341700e+04	E
24	Argentina	9.0	Population density	4107.0	1962.0	7.656000e+00	E
25	Argentina	9.0	Population density	4107.0	1967.0	8.248000e+00	E
26	Argentina	9.0	Population density	4107.0	1972.0	8.913000e+00	E
27	Argentina	9.0	Population density	4107.0	1977.0	9.667000e+00	E
28	Argentina	9.0	Population density	4107.0	1982.0	1.043000e+01	E
29	Argentina	9.0	Population density	4107.0	1987.0	1.127000e+01	E
...	...	...	...	...	...	...	...



	Area	Area Id	Variable Name	Variable Id	Year	Value	Symbol
360	United States of America	231.0	Population density	4107.0	1972.0	2.214000e+01	E
361	United States of America	231.0	Population density	4107.0	1977.0	2.317000e+01	E
362	United States of America	231.0	Population density	4107.0	1982.0	2.430000e+01	E
363	United States of America	231.0	Population density	4107.0	1987.0	2.549000e+01	E
364	United States of America	231.0	Population density	4107.0	1992.0	2.678000e+01	E
365	United States of America	231.0	Population density	4107.0	1997.0	2.834000e+01	E
366	United States of America	231.0	Population density	4107.0	2002.0	2.995000e+01	E
367	United States of America	231.0	Population density	4107.0	2007.0	3.132000e+01	E
368	United States of America	231.0	Population density	4107.0	2012.0	3.202000e+01	E
369	United States of America	231.0	Population density	4107.0	2015.0	3.273000e+01	E
370	United States of America	231.0	Gross Domestic Product (GDP)	4112.0	1962.0	6.050000e+11	E
371	United States of America	231.0	Gross Domestic Product (GDP)	4112.0	1967.0	8.620000e+11	E
372	United States of America	231.0	Gross Domestic Product (GDP)	4112.0	1972.0	1.280000e+12	E
373	United States of America	231.0	Gross Domestic Product (GDP)	4112.0	1977.0	2.090000e+12	E
374	United States of America	231.0	Gross Domestic Product (GDP)	4112.0	1982.0	3.340000e+12	E
375	United States of America	231.0	Gross Domestic Product (GDP)	4112.0	1987.0	4.870000e+12	E
376	United States of America	231.0	Gross Domestic Product (GDP)	4112.0	1992.0	6.540000e+12	E
377	United States of America	231.0	Gross Domestic Product (GDP)	4112.0	1997.0	8.610000e+12	E
378	United States of America	231.0	Gross Domestic Product (GDP)	4112.0	2002.0	1.100000e+13	E
379	United States of America	231.0	Gross Domestic Product (GDP)	4112.0	2007.0	1.450000e+13	E
380	United States of America	231.0	Gross Domestic Product (GDP)	4112.0	2012.0	1.620000e+13	E
381	United States of America	231.0	Gross Domestic Product (GDP)	4112.0	2015.0	1.790000e+13	E
382	United States of America	231.0	National Rainfall Index (NRI)	4472.0	1965.0	9.285000e+02	E
383	United States of America	231.0	National Rainfall Index (NRI)	4472.0	1969.0	9.522000e+02	E
384	United States of America	231.0	National Rainfall Index (NRI)	4472.0	1974.0	1.008000e+03	E

	Area	Area Id	Variable Name	Variable Id	Year	Value	Symbol
385	United States of America	231.0	National Rainfall Index (NRI)	4472.0	1981.0	9.492000e+02	E
386	United States of America	231.0	National Rainfall Index (NRI)	4472.0	1984.0	9.746000e+02	E
387	United States of America	231.0	National Rainfall Index (NRI)	4472.0	1992.0	1.020000e+03	E
388	United States of America	231.0	National Rainfall Index (NRI)	4472.0	1996.0	1.005000e+03	E
389	United States of America	231.0	National Rainfall Index (NRI)	4472.0	2002.0	9.387000e+02	E

390 rows × 7 columns

**6a) For our analysis we do not want all the columns in our dataframe. Lets drop all the redundant columns/ features.**

**Drop columns: Area Id, Variable Id, Symbol. Save the new dataframe as df1. Display the first 5 rows of the new dataframe.**

In [28]:

```
# your code here
data_no_nan = data_no_nan.drop('Area Id', axis = 1)
data_no_nan = data_no_nan.drop('Variable Id', axis = 1)
data_no_nan = data_no_nan.drop('Symbol', axis = 1)
df1 = data_no_nan
df1
```

Out[28]:

	Area	Variable Name	Year	Value
0	Argentina	Total area of the country	1962.0	2.780400e+05
1	Argentina	Total area of the country	1967.0	2.780400e+05
2	Argentina	Total area of the country	1972.0	2.780400e+05
3	Argentina	Total area of the country	1977.0	2.780400e+05
4	Argentina	Total area of the country	1982.0	2.780400e+05
5	Argentina	Total area of the country	1987.0	2.780400e+05
6	Argentina	Total area of the country	1992.0	2.780400e+05
7	Argentina	Total area of the country	1997.0	2.780400e+05
8	Argentina	Total area of the country	2002.0	2.780400e+05
9	Argentina	Total area of the country	2007.0	2.780400e+05
10	Argentina	Total area of the country	2012.0	2.780400e+05
11	Argentina	Total area of the country	2014.0	2.780400e+05
12	Argentina	Total population	1962.0	2.128800e+04
13	Argentina	Total population	1967.0	2.293200e+04
14	Argentina	Total population	1972.0	2.478300e+04
15	Argentina	Total population	1977.0	2.687900e+04
16	Argentina	Total population	1982.0	2.899400e+04
17	Argentina	Total population	1987.0	3.132600e+04
18	Argentina	Total population	1992.0	3.365500e+04
19	Argentina	Total population	1997.0	3.583400e+04
20	Argentina	Total population	2002.0	3.788900e+04
21	Argentina	Total population	2007.0	3.997000e+04
22	Argentina	Total population	2012.0	4.209500e+04
23	Argentina	Total population	2015.0	4.341700e+04
24	Argentina	Population density	1962.0	7.656000e+00
25	Argentina	Population density	1967.0	8.248000e+00
26	Argentina	Population density	1972.0	8.913000e+00
27	Argentina	Population density	1977.0	9.667000e+00
28	Argentina	Population density	1982.0	1.043000e+01
29	Argentina	Population density	1987.0	1.127000e+01

	Area	Variable Name	Year	Value
...	...	...	...	...
360	United States of America	Population density	1972.0	2.214000e+01
361	United States of America	Population density	1977.0	2.317000e+01
362	United States of America	Population density	1982.0	2.430000e+01
363	United States of America	Population density	1987.0	2.549000e+01
364	United States of America	Population density	1992.0	2.678000e+01
365	United States of America	Population density	1997.0	2.834000e+01
366	United States of America	Population density	2002.0	2.995000e+01
367	United States of America	Population density	2007.0	3.132000e+01
368	United States of America	Population density	2012.0	3.202000e+01
369	United States of America	Population density	2015.0	3.273000e+01
370	United States of America	Gross Domestic Product (GDP)	1962.0	6.050000e+11
371	United States of America	Gross Domestic Product (GDP)	1967.0	8.620000e+11
372	United States of America	Gross Domestic Product (GDP)	1972.0	1.280000e+12
373	United States of America	Gross Domestic Product (GDP)	1977.0	2.090000e+12
374	United States of America	Gross Domestic Product (GDP)	1982.0	3.340000e+12
375	United States of America	Gross Domestic Product (GDP)	1987.0	4.870000e+12
376	United States of America	Gross Domestic Product (GDP)	1992.0	6.540000e+12
377	United States of America	Gross Domestic Product (GDP)	1997.0	8.610000e+12
378	United States of America	Gross Domestic Product (GDP)	2002.0	1.100000e+13
379	United States of America	Gross Domestic Product (GDP)	2007.0	1.450000e+13
380	United States of America	Gross Domestic Product (GDP)	2012.0	1.620000e+13
381	United States of America	Gross Domestic Product (GDP)	2015.0	1.790000e+13
382	United States of America	National Rainfall Index (NRI)	1965.0	9.285000e+02
383	United States of America	National Rainfall Index (NRI)	1969.0	9.522000e+02
384	United States of America	National Rainfall Index (NRI)	1974.0	1.008000e+03
385	United States of America	National Rainfall Index (NRI)	1981.0	9.492000e+02
386	United States of America	National Rainfall Index (NRI)	1984.0	9.746000e+02
387	United States of America	National Rainfall Index (NRI)	1992.0	1.020000e+03
388	United States of America	National Rainfall Index (NRI)	1996.0	1.005000e+03
389	United States of America	National Rainfall Index (NRI)	2002.0	9.387000e+02

390 rows × 4 columns

**6b) Display all the unique values in your new dataframe for columns: Area, Variable Name, Year.**

In [29]:

```
# your code here
df1['Area'].unique()
```

Out[29]:

```
array(['Argentina', 'Australia', 'Germany', 'Iceland', 'Ireland',
      'Sweden', 'United States of America'], dtype=object)
```

In [30]:

```
df1['Variable Name'].unique()
```

Out[30]:

```
array(['Total area of the country', 'Total population',
      'Population density', 'Gross Domestic Product (GDP)',
      'National Rainfall Index (NRI)'], dtype=object)
```

In [31]:

```
df1['Year'].unique()
```

Out[31]:

```
array([1962., 1967., 1972., 1977., 1982., 1987., 1992., 1997., 2002.,
      2007., 2012., 2014., 2015., 1963., 1970., 1974., 1978., 1984.,
      1990., 1964., 1981., 1985., 1996., 2001., 1969., 1973., 1979.,
      1993., 1971., 1975., 1986., 1991., 1998., 2000., 1965., 1983.,
      1988., 1995.])
```

**6c) Convert the year column to pandas datetime. Convert the 'Year' column float values to pandas datetime objects, where each year is represented as the first day of that year. Also display the column and datatype for 'Year' after conversion. For eg: 1962.0 will be represented as 1962-01-01**

In [32]:

```
df1['Year'] = pd.to_datetime(df1['Year'], format = "%Y.0")
s1 = df1['Year']
print(type(s1[1]))
df1
```

```
<class 'pandas._libs.tslibs.timestamps.Timestamp'>
```

Out[32]:

	Area	Variable Name	Year	Value
0	Argentina	Total area of the country	1962-01-01	2.780400e+05
1	Argentina	Total area of the country	1967-01-01	2.780400e+05
2	Argentina	Total area of the country	1972-01-01	2.780400e+05
3	Argentina	Total area of the country	1977-01-01	2.780400e+05
4	Argentina	Total area of the country	1982-01-01	2.780400e+05
5	Argentina	Total area of the country	1987-01-01	2.780400e+05
6	Argentina	Total area of the country	1992-01-01	2.780400e+05
7	Argentina	Total area of the country	1997-01-01	2.780400e+05
8	Argentina	Total area of the country	2002-01-01	2.780400e+05
9	Argentina	Total area of the country	2007-01-01	2.780400e+05
10	Argentina	Total area of the country	2012-01-01	2.780400e+05
11	Argentina	Total area of the country	2014-01-01	2.780400e+05
12	Argentina	Total population	1962-01-01	2.128800e+04
13	Argentina	Total population	1967-01-01	2.293200e+04
14	Argentina	Total population	1972-01-01	2.478300e+04
15	Argentina	Total population	1977-01-01	2.687900e+04
16	Argentina	Total population	1982-01-01	2.899400e+04
17	Argentina	Total population	1987-01-01	3.132600e+04
18	Argentina	Total population	1992-01-01	3.365500e+04
19	Argentina	Total population	1997-01-01	3.583400e+04
20	Argentina	Total population	2002-01-01	3.788900e+04
21	Argentina	Total population	2007-01-01	3.997000e+04
22	Argentina	Total population	2012-01-01	4.209500e+04
23	Argentina	Total population	2015-01-01	4.341700e+04
24	Argentina	Population density	1962-01-01	7.656000e+00
25	Argentina	Population density	1967-01-01	8.248000e+00
26	Argentina	Population density	1972-01-01	8.913000e+00
27	Argentina	Population density	1977-01-01	9.667000e+00
28	Argentina	Population density	1982-01-01	1.043000e+01
29	Argentina	Population density	1987-01-01	1.127000e+01

	Area	Variable Name	Year	Value
...	...	...	...	...
360	United States of America	Population density	1972-01-01	2.214000e+01
361	United States of America	Population density	1977-01-01	2.317000e+01
362	United States of America	Population density	1982-01-01	2.430000e+01
363	United States of America	Population density	1987-01-01	2.549000e+01
364	United States of America	Population density	1992-01-01	2.678000e+01
365	United States of America	Population density	1997-01-01	2.834000e+01
366	United States of America	Population density	2002-01-01	2.995000e+01
367	United States of America	Population density	2007-01-01	3.132000e+01
368	United States of America	Population density	2012-01-01	3.202000e+01
369	United States of America	Population density	2015-01-01	3.273000e+01
370	United States of America	Gross Domestic Product (GDP)	1962-01-01	6.050000e+11
371	United States of America	Gross Domestic Product (GDP)	1967-01-01	8.620000e+11
372	United States of America	Gross Domestic Product (GDP)	1972-01-01	1.280000e+12
373	United States of America	Gross Domestic Product (GDP)	1977-01-01	2.090000e+12
374	United States of America	Gross Domestic Product (GDP)	1982-01-01	3.340000e+12
375	United States of America	Gross Domestic Product (GDP)	1987-01-01	4.870000e+12
376	United States of America	Gross Domestic Product (GDP)	1992-01-01	6.540000e+12
377	United States of America	Gross Domestic Product (GDP)	1997-01-01	8.610000e+12
378	United States of America	Gross Domestic Product (GDP)	2002-01-01	1.100000e+13
379	United States of America	Gross Domestic Product (GDP)	2007-01-01	1.450000e+13
380	United States of America	Gross Domestic Product (GDP)	2012-01-01	1.620000e+13
381	United States of America	Gross Domestic Product (GDP)	2015-01-01	1.790000e+13
382	United States of America	National Rainfall Index (NRI)	1965-01-01	9.285000e+02
383	United States of America	National Rainfall Index (NRI)	1969-01-01	9.522000e+02
384	United States of America	National Rainfall Index (NRI)	1974-01-01	1.008000e+03
385	United States of America	National Rainfall Index (NRI)	1981-01-01	9.492000e+02
386	United States of America	National Rainfall Index (NRI)	1984-01-01	9.746000e+02
387	United States of America	National Rainfall Index (NRI)	1992-01-01	1.020000e+03
388	United States of America	National Rainfall Index (NRI)	1996-01-01	1.005000e+03
389	United States of America	National Rainfall Index (NRI)	2002-01-01	9.387000e+02

390 rows × 4 columns

## Extract specific statistics from the preprocessed data:

7a) Create a dataframe 'dftemp' to store rows where Area is 'Iceland'. Display the dataframe.

In [33]:

```
# your code here
dftemp = df1[:, df1['Area'] == 'Iceland']
dftemp
```

Out[33]:

	Area	Variable Name	Year	Value
166	Iceland	Total area of the country	1962-01-01	1.030000e+04
167	Iceland	Total area of the country	1967-01-01	1.030000e+04
168	Iceland	Total area of the country	1972-01-01	1.030000e+04
169	Iceland	Total area of the country	1977-01-01	1.030000e+04
170	Iceland	Total area of the country	1982-01-01	1.030000e+04
171	Iceland	Total area of the country	1987-01-01	1.030000e+04
172	Iceland	Total area of the country	1992-01-01	1.030000e+04
173	Iceland	Total area of the country	1997-01-01	1.030000e+04
174	Iceland	Total area of the country	2002-01-01	1.030000e+04
175	Iceland	Total area of the country	2007-01-01	1.030000e+04
176	Iceland	Total area of the country	2012-01-01	1.030000e+04
177	Iceland	Total area of the country	2014-01-01	1.030000e+04
178	Iceland	Total population	1962-01-01	1.826000e+02
179	Iceland	Total population	1967-01-01	1.974000e+02
180	Iceland	Total population	1972-01-01	2.099000e+02
181	Iceland	Total population	1977-01-01	2.221000e+02
182	Iceland	Total population	1982-01-01	2.331000e+02
183	Iceland	Total population	1987-01-01	2.469000e+02
184	Iceland	Total population	1992-01-01	2.599000e+02
185	Iceland	Total population	1997-01-01	2.728000e+02
186	Iceland	Total population	2002-01-01	2.869000e+02
187	Iceland	Total population	2007-01-01	3.054000e+02
188	Iceland	Total population	2012-01-01	3.234000e+02
189	Iceland	Total population	2015-01-01	3.294000e+02
190	Iceland	Population density	1962-01-01	1.773000e+00
191	Iceland	Population density	1967-01-01	1.917000e+00
192	Iceland	Population density	1972-01-01	2.038000e+00
193	Iceland	Population density	1977-01-01	2.156000e+00
194	Iceland	Population density	1982-01-01	2.263000e+00
195	Iceland	Population density	1987-01-01	2.397000e+00
196	Iceland	Population density	1992-01-01	2.523000e+00
197	Iceland	Population density	1997-01-01	2.649000e+00



	Area	Variable Name	Year	Value
198	Iceland	Population density	2002-01-01	2.785000e+00
199	Iceland	Population density	2007-01-01	2.965000e+00
200	Iceland	Population density	2012-01-01	3.140000e+00
201	Iceland	Population density	2015-01-01	3.198000e+00
202	Iceland	Gross Domestic Product (GDP)	1962-01-01	2.849165e+08
203	Iceland	Gross Domestic Product (GDP)	1967-01-01	6.212260e+08
204	Iceland	Gross Domestic Product (GDP)	1972-01-01	8.465069e+08
205	Iceland	Gross Domestic Product (GDP)	1977-01-01	2.226539e+09
206	Iceland	Gross Domestic Product (GDP)	1982-01-01	3.232804e+09
207	Iceland	Gross Domestic Product (GDP)	1987-01-01	5.565384e+09
208	Iceland	Gross Domestic Product (GDP)	1992-01-01	7.138788e+09
209	Iceland	Gross Domestic Product (GDP)	1997-01-01	7.596126e+09
210	Iceland	Gross Domestic Product (GDP)	2002-01-01	9.161798e+09
211	Iceland	Gross Domestic Product (GDP)	2007-01-01	2.129384e+10
212	Iceland	Gross Domestic Product (GDP)	2012-01-01	1.419452e+10
213	Iceland	Gross Domestic Product (GDP)	2015-01-01	1.659849e+10
214	Iceland	National Rainfall Index (NRI)	1967-01-01	8.160000e+02
215	Iceland	National Rainfall Index (NRI)	1971-01-01	9.632000e+02
216	Iceland	National Rainfall Index (NRI)	1975-01-01	1.010000e+03
217	Iceland	National Rainfall Index (NRI)	1981-01-01	9.326000e+02
218	Iceland	National Rainfall Index (NRI)	1986-01-01	9.685000e+02
219	Iceland	National Rainfall Index (NRI)	1991-01-01	1.095000e+03
220	Iceland	National Rainfall Index (NRI)	1997-01-01	9.932000e+02
221	Iceland	National Rainfall Index (NRI)	1998-01-01	9.234000e+02

**7b) Print the years when the National Rainfall Index (NRI) was greater than 900 and less than 950 in Iceland. Use the dataframe you created in the previous question 'dftemp'.**

In [34]:

```
# your code here
dftemp = dftemp[dftemp['Variable Name'] == 'National Rainfall Index (NRI)']
dftemp = dftemp[dftemp['Value'] > 900]
dftemp = dftemp[dftemp['Value'] < 950]
dftemp
```

Out[34]:

	Area	Variable Name	Year	Value
217	Iceland	National Rainfall Index (NRI)	1981-01-01	932.6
221	Iceland	National Rainfall Index (NRI)	1998-01-01	923.4

In [35]:

```
print("years when the National Rainfall Index (NRI) was greater than 900 and less than 950 in Iceland are:")  
dftemp['Year']
```

years when the National Rainfall Index (NRI) was greater than 900 and less than 950 in Iceland are:

Out[35]:

```
217    1981-01-01  
221    1998-01-01  
Name: Year, dtype: datetime64[ns]
```

## US statistics:

8a) Create a new DataFrame called `df_usa` that only contains values where 'Area' is equal to 'United States of America'. Set the indices to be the 'Year' column ( Use `.set_index()` ). Display the dataframe head.

In [36]:

```
# your code here
df_usa = df1[:, df1['Area'] == 'United States of America']
df_usa = df_usa.set_index('Year')
df_usa
```

Out[36]:

	Area	Variable Name	Value
Year			
1962-01-01	United States of America	Total area of the country	9.629090e+05
1967-01-01	United States of America	Total area of the country	9.629090e+05
1972-01-01	United States of America	Total area of the country	9.629090e+05
1977-01-01	United States of America	Total area of the country	9.629090e+05
1982-01-01	United States of America	Total area of the country	9.629090e+05
1987-01-01	United States of America	Total area of the country	9.629090e+05
1992-01-01	United States of America	Total area of the country	9.629090e+05
1997-01-01	United States of America	Total area of the country	9.629090e+05
2002-01-01	United States of America	Total area of the country	9.632030e+05
2007-01-01	United States of America	Total area of the country	9.632030e+05
2012-01-01	United States of America	Total area of the country	9.831510e+05
2014-01-01	United States of America	Total area of the country	9.831510e+05
1962-01-01	United States of America	Total population	1.918610e+05
1967-01-01	United States of America	Total population	2.037130e+05
1972-01-01	United States of America	Total population	2.132200e+05
1977-01-01	United States of America	Total population	2.230910e+05
1982-01-01	United States of America	Total population	2.339540e+05
1987-01-01	United States of America	Total population	2.454250e+05
1992-01-01	United States of America	Total population	2.579080e+05
1997-01-01	United States of America	Total population	2.728830e+05
2002-01-01	United States of America	Total population	2.884710e+05
2007-01-01	United States of America	Total population	3.016560e+05
2012-01-01	United States of America	Total population	3.147990e+05
2015-01-01	United States of America	Total population	3.217740e+05
1962-01-01	United States of America	Population density	1.993000e+01
1967-01-01	United States of America	Population density	2.116000e+01
1972-01-01	United States of America	Population density	2.214000e+01
1977-01-01	United States of America	Population density	2.317000e+01
1982-01-01	United States of America	Population density	2.430000e+01
1987-01-01	United States of America	Population density	2.549000e+01

	Area	Variable Name	Value
Year			
1992-01-01	United States of America	Population density	2.678000e+01
1997-01-01	United States of America	Population density	2.834000e+01
2002-01-01	United States of America	Population density	2.995000e+01
2007-01-01	United States of America	Population density	3.132000e+01
2012-01-01	United States of America	Population density	3.202000e+01
2015-01-01	United States of America	Population density	3.273000e+01
1962-01-01	United States of America	Gross Domestic Product (GDP)	6.050000e+11
1967-01-01	United States of America	Gross Domestic Product (GDP)	8.620000e+11
1972-01-01	United States of America	Gross Domestic Product (GDP)	1.280000e+12
1977-01-01	United States of America	Gross Domestic Product (GDP)	2.090000e+12
1982-01-01	United States of America	Gross Domestic Product (GDP)	3.340000e+12
1987-01-01	United States of America	Gross Domestic Product (GDP)	4.870000e+12
1992-01-01	United States of America	Gross Domestic Product (GDP)	6.540000e+12
1997-01-01	United States of America	Gross Domestic Product (GDP)	8.610000e+12
2002-01-01	United States of America	Gross Domestic Product (GDP)	1.100000e+13
2007-01-01	United States of America	Gross Domestic Product (GDP)	1.450000e+13
2012-01-01	United States of America	Gross Domestic Product (GDP)	1.620000e+13
2015-01-01	United States of America	Gross Domestic Product (GDP)	1.790000e+13
1965-01-01	United States of America	National Rainfall Index (NRI)	9.285000e+02
1969-01-01	United States of America	National Rainfall Index (NRI)	9.522000e+02
1974-01-01	United States of America	National Rainfall Index (NRI)	1.008000e+03
1981-01-01	United States of America	National Rainfall Index (NRI)	9.492000e+02
1984-01-01	United States of America	National Rainfall Index (NRI)	9.746000e+02
1992-01-01	United States of America	National Rainfall Index (NRI)	1.020000e+03
1996-01-01	United States of America	National Rainfall Index (NRI)	1.005000e+03
2002-01-01	United States of America	National Rainfall Index (NRI)	9.387000e+02

**8b) Pivot the DataFrame so that the unique values in the column 'Variable Name' becomes the columns. The DataFrame values should be the ones in the the 'Value' column. Save it in df\_usa. Display the dataframe head.**

In [37]:

```
# your code here
new_col = df_usa['Variable Name'].unique()
for i in new_col:
    df_usa[i] = df_usa['Value'][df_usa["Variable Name"] == i]

df_usa.head()
```

Out[37]:

	Area	Variable Name	Value	Total area of the country	Total population	Population density	Gross Domestic Product (GDP)	National Rainfall Index (NRI)
Year								
1962-01-01	United States of America	Total area of the country	962909.0	962909.0	191861.0	19.93	6.050000e+11	NaN
1967-01-01	United States of America	Total area of the country	962909.0	962909.0	203713.0	21.16	8.620000e+11	NaN
1972-01-01	United States of America	Total area of the country	962909.0	962909.0	213220.0	22.14	1.280000e+12	NaN
1977-01-01	United States of America	Total area of the country	962909.0	962909.0	223091.0	23.17	2.090000e+12	NaN
1982-01-01	United States of America	Total area of the country	962909.0	962909.0	233954.0	24.30	3.340000e+12	NaN

8c) Rename new columns to ['GDP','NRI','PD','Area','Population'] and display the head.

In [38]:

```
# your code her
df_usa.columns = [ 'Area' , 'Variable Name' , 'Value' , 'Phy Area' , 'Population' , 'Popula
df_usa = df_usa.drop(columns = [ 'Variable Name' ] , axis = 1)
df_usa = df_usa.drop_duplicates(subset= None, keep=False, inplace=False)
df_usa.head()
```

Out[38]:

	Area	Value	Phy Area	Population	Population density	GDP	NRI
Year							
1962-01-01	United States of America	962909.0	962909.0	191861.0	19.93	6.050000e+11	NaN
1967-01-01	United States of America	962909.0	962909.0	203713.0	21.16	8.620000e+11	NaN
1972-01-01	United States of America	962909.0	962909.0	213220.0	22.14	1.280000e+12	NaN
1977-01-01	United States of America	962909.0	962909.0	223091.0	23.17	2.090000e+12	NaN
1982-01-01	United States of America	962909.0	962909.0	233954.0	24.30	3.340000e+12	NaN

8d) Replace all 'Nan' values in df\_usa with 0. Display the head of the dataframe.

In [39]:

```
# your code here
df_usa= df_usa.fillna(0)
df_usa
```

Out[39]:

	Area	Value	Phy Area	Population	Population density	GDP	NRI
Year							
1962-01-01	United States of America	9.629090e+05	962909.0	191861.0	19.93	6.050000e+11	0.0
1967-01-01	United States of America	9.629090e+05	962909.0	203713.0	21.16	8.620000e+11	0.0
1972-01-01	United States of America	9.629090e+05	962909.0	213220.0	22.14	1.280000e+12	0.0
1977-01-01	United States of America	9.629090e+05	962909.0	223091.0	23.17	2.090000e+12	0.0
1982-01-01	United States of America	9.629090e+05	962909.0	233954.0	24.30	3.340000e+12	0.0
1987-01-01	United States of America	9.629090e+05	962909.0	245425.0	25.49	4.870000e+12	0.0
1992-01-01	United States of America	9.629090e+05	962909.0	257908.0	26.78	6.540000e+12	1020.0
1997-01-01	United States of America	9.629090e+05	962909.0	272883.0	28.34	8.610000e+12	0.0
2002-01-01	United States of America	9.632030e+05	963203.0	288471.0	29.95	1.100000e+13	938.7
2007-01-01	United States of America	9.632030e+05	963203.0	301656.0	31.32	1.450000e+13	0.0
2012-01-01	United States of America	9.831510e+05	983151.0	314799.0	32.02	1.620000e+13	0.0
2014-01-01	United States of America	9.831510e+05	983151.0	0.0	0.00	0.000000e+00	0.0
1962-01-01	United States of America	1.918610e+05	962909.0	191861.0	19.93	6.050000e+11	0.0
1967-01-01	United States of America	2.037130e+05	962909.0	203713.0	21.16	8.620000e+11	0.0
1972-01-01	United States of America	2.132200e+05	962909.0	213220.0	22.14	1.280000e+12	0.0
1977-01-01	United States of America	2.230910e+05	962909.0	223091.0	23.17	2.090000e+12	0.0
1982-01-01	United States of America	2.339540e+05	962909.0	233954.0	24.30	3.340000e+12	0.0
1987-01-01	United States of America	2.454250e+05	962909.0	245425.0	25.49	4.870000e+12	0.0
1992-01-01	United States of America	2.579080e+05	962909.0	257908.0	26.78	6.540000e+12	1020.0

	Area	Value	Phy Area	Population	Population density	GDP	NRI
Year							
1997-01-01	United States of America	2.728830e+05	962909.0	272883.0	28.34	8.610000e+12	0.0
2002-01-01	United States of America	2.884710e+05	963203.0	288471.0	29.95	1.100000e+13	938.7
2007-01-01	United States of America	3.016560e+05	963203.0	301656.0	31.32	1.450000e+13	0.0
2012-01-01	United States of America	3.147990e+05	983151.0	314799.0	32.02	1.620000e+13	0.0
2015-01-01	United States of America	3.217740e+05	0.0	321774.0	32.73	1.790000e+13	0.0
1962-01-01	United States of America	1.993000e+01	962909.0	191861.0	19.93	6.050000e+11	0.0
1967-01-01	United States of America	2.116000e+01	962909.0	203713.0	21.16	8.620000e+11	0.0
1972-01-01	United States of America	2.214000e+01	962909.0	213220.0	22.14	1.280000e+12	0.0
1977-01-01	United States of America	2.317000e+01	962909.0	223091.0	23.17	2.090000e+12	0.0
1982-01-01	United States of America	2.430000e+01	962909.0	233954.0	24.30	3.340000e+12	0.0
1987-01-01	United States of America	2.549000e+01	962909.0	245425.0	25.49	4.870000e+12	0.0
1992-01-01	United States of America	2.678000e+01	962909.0	257908.0	26.78	6.540000e+12	1020.0
1997-01-01	United States of America	2.834000e+01	962909.0	272883.0	28.34	8.610000e+12	0.0
2002-01-01	United States of America	2.995000e+01	963203.0	288471.0	29.95	1.100000e+13	938.7
2007-01-01	United States of America	3.132000e+01	963203.0	301656.0	31.32	1.450000e+13	0.0
2012-01-01	United States of America	3.202000e+01	983151.0	314799.0	32.02	1.620000e+13	0.0
2015-01-01	United States of America	3.273000e+01	0.0	321774.0	32.73	1.790000e+13	0.0
1962-01-01	United States of America	6.050000e+11	962909.0	191861.0	19.93	6.050000e+11	0.0
1967-01-01	United States of America	8.620000e+11	962909.0	203713.0	21.16	8.620000e+11	0.0
1972-01-01	United States of America	1.280000e+12	962909.0	213220.0	22.14	1.280000e+12	0.0
1977-01-01	United States of America	2.090000e+12	962909.0	223091.0	23.17	2.090000e+12	0.0
1982-01-01	United States of America	3.340000e+12	962909.0	233954.0	24.30	3.340000e+12	0.0
1987-01-01	United States of America	4.870000e+12	962909.0	245425.0	25.49	4.870000e+12	0.0



	Area	Value	Phy Area	Population	Population density	GDP	NRI
Year							
1992-01-01	United States of America	6.540000e+12	962909.0	257908.0	26.78	6.540000e+12	1020.0
1997-01-01	United States of America	8.610000e+12	962909.0	272883.0	28.34	8.610000e+12	0.0
2002-01-01	United States of America	1.100000e+13	963203.0	288471.0	29.95	1.100000e+13	938.7
2007-01-01	United States of America	1.450000e+13	963203.0	301656.0	31.32	1.450000e+13	0.0
2012-01-01	United States of America	1.620000e+13	983151.0	314799.0	32.02	1.620000e+13	0.0
2015-01-01	United States of America	1.790000e+13	0.0	321774.0	32.73	1.790000e+13	0.0
1965-01-01	United States of America	9.285000e+02	0.0	0.0	0.00	0.000000e+00	928.5
1969-01-01	United States of America	9.522000e+02	0.0	0.0	0.00	0.000000e+00	952.2
1974-01-01	United States of America	1.008000e+03	0.0	0.0	0.00	0.000000e+00	1008.0
1981-01-01	United States of America	9.492000e+02	0.0	0.0	0.00	0.000000e+00	949.2
1984-01-01	United States of America	9.746000e+02	0.0	0.0	0.00	0.000000e+00	974.6
1992-01-01	United States of America	1.020000e+03	962909.0	257908.0	26.78	6.540000e+12	1020.0
1996-01-01	United States of America	1.005000e+03	0.0	0.0	0.00	0.000000e+00	1005.0
2002-01-01	United States of America	9.387000e+02	963203.0	288471.0	29.95	1.100000e+13	938.7

## Note: Use df\_usa

9a) Multiply the 'Area' column for all countries by 10 (so instead of 1000 ha, the unit becomes 100 ha = 1km<sup>2</sup>). Display the dataframe head.

In [40]:

```
# your code here
df_usa[ 'Phy Area' ] = df_usa[ 'Phy Area' ]*10
df_usa
```

Out[40]:

	Area	Value	Phy Area	Population	Population density	GDP	NRI
Year							
1962-01-01	United States of America	9.629090e+05	9629090.0	191861.0	19.93	6.050000e+11	0.0
1967-01-01	United States of America	9.629090e+05	9629090.0	203713.0	21.16	8.620000e+11	0.0
1972-01-01	United States of America	9.629090e+05	9629090.0	213220.0	22.14	1.280000e+12	0.0
1977-01-01	United States of America	9.629090e+05	9629090.0	223091.0	23.17	2.090000e+12	0.0
1982-01-01	United States of America	9.629090e+05	9629090.0	233954.0	24.30	3.340000e+12	0.0
1987-01-01	United States of America	9.629090e+05	9629090.0	245425.0	25.49	4.870000e+12	0.0
1992-01-01	United States of America	9.629090e+05	9629090.0	257908.0	26.78	6.540000e+12	1020.0
1997-01-01	United States of America	9.629090e+05	9629090.0	272883.0	28.34	8.610000e+12	0.0
2002-01-01	United States of America	9.632030e+05	9632030.0	288471.0	29.95	1.100000e+13	938.7
2007-01-01	United States of America	9.632030e+05	9632030.0	301656.0	31.32	1.450000e+13	0.0
2012-01-01	United States of America	9.831510e+05	9831510.0	314799.0	32.02	1.620000e+13	0.0
2014-01-01	United States of America	9.831510e+05	9831510.0	0.0	0.00	0.000000e+00	0.0
1962-01-01	United States of America	1.918610e+05	9629090.0	191861.0	19.93	6.050000e+11	0.0
1967-01-01	United States of America	2.037130e+05	9629090.0	203713.0	21.16	8.620000e+11	0.0
1972-01-01	United States of America	2.132200e+05	9629090.0	213220.0	22.14	1.280000e+12	0.0
1977-01-01	United States of America	2.230910e+05	9629090.0	223091.0	23.17	2.090000e+12	0.0
1982-01-01	United States of America	2.339540e+05	9629090.0	233954.0	24.30	3.340000e+12	0.0
1987-01-01	United States of America	2.454250e+05	9629090.0	245425.0	25.49	4.870000e+12	0.0
1992-01-01	United States of America	2.579080e+05	9629090.0	257908.0	26.78	6.540000e+12	1020.0

	Area	Value	Phy Area	Population	Population density	GDP	NRI
Year							
1997-01-01	United States of America	2.728830e+05	9629090.0	272883.0	28.34	8.610000e+12	0.0
2002-01-01	United States of America	2.884710e+05	9632030.0	288471.0	29.95	1.100000e+13	938.7
2007-01-01	United States of America	3.016560e+05	9632030.0	301656.0	31.32	1.450000e+13	0.0
2012-01-01	United States of America	3.147990e+05	9831510.0	314799.0	32.02	1.620000e+13	0.0
2015-01-01	United States of America	3.217740e+05	0.0	321774.0	32.73	1.790000e+13	0.0
1962-01-01	United States of America	1.993000e+01	9629090.0	191861.0	19.93	6.050000e+11	0.0
1967-01-01	United States of America	2.116000e+01	9629090.0	203713.0	21.16	8.620000e+11	0.0
1972-01-01	United States of America	2.214000e+01	9629090.0	213220.0	22.14	1.280000e+12	0.0
1977-01-01	United States of America	2.317000e+01	9629090.0	223091.0	23.17	2.090000e+12	0.0
1982-01-01	United States of America	2.430000e+01	9629090.0	233954.0	24.30	3.340000e+12	0.0
1987-01-01	United States of America	2.549000e+01	9629090.0	245425.0	25.49	4.870000e+12	0.0
1992-01-01	United States of America	2.678000e+01	9629090.0	257908.0	26.78	6.540000e+12	1020.0
1997-01-01	United States of America	2.834000e+01	9629090.0	272883.0	28.34	8.610000e+12	0.0
2002-01-01	United States of America	2.995000e+01	9632030.0	288471.0	29.95	1.100000e+13	938.7
2007-01-01	United States of America	3.132000e+01	9632030.0	301656.0	31.32	1.450000e+13	0.0
2012-01-01	United States of America	3.202000e+01	9831510.0	314799.0	32.02	1.620000e+13	0.0
2015-01-01	United States of America	3.273000e+01	0.0	321774.0	32.73	1.790000e+13	0.0
1962-01-01	United States of America	6.050000e+11	9629090.0	191861.0	19.93	6.050000e+11	0.0
1967-01-01	United States of America	8.620000e+11	9629090.0	203713.0	21.16	8.620000e+11	0.0
1972-01-01	United States of America	1.280000e+12	9629090.0	213220.0	22.14	1.280000e+12	0.0
1977-01-01	United States of America	2.090000e+12	9629090.0	223091.0	23.17	2.090000e+12	0.0
1982-01-01	United States of America	3.340000e+12	9629090.0	233954.0	24.30	3.340000e+12	0.0
1987-01-01	United States of America	4.870000e+12	9629090.0	245425.0	25.49	4.870000e+12	0.0

	Area	Value	Phy Area	Population	Population density	GDP	NRI
Year							
1992-01-01	United States of America	6.540000e+12	9629090.0	257908.0	26.78	6.540000e+12	1020.0
1997-01-01	United States of America	8.610000e+12	9629090.0	272883.0	28.34	8.610000e+12	0.0
2002-01-01	United States of America	1.100000e+13	9632030.0	288471.0	29.95	1.100000e+13	938.7
2007-01-01	United States of America	1.450000e+13	9632030.0	301656.0	31.32	1.450000e+13	0.0
2012-01-01	United States of America	1.620000e+13	9831510.0	314799.0	32.02	1.620000e+13	0.0
2015-01-01	United States of America	1.790000e+13	0.0	321774.0	32.73	1.790000e+13	0.0
1965-01-01	United States of America	9.285000e+02	0.0	0.0	0.00	0.000000e+00	928.5
1969-01-01	United States of America	9.522000e+02	0.0	0.0	0.00	0.000000e+00	952.2
1974-01-01	United States of America	1.008000e+03	0.0	0.0	0.00	0.000000e+00	1008.0
1981-01-01	United States of America	9.492000e+02	0.0	0.0	0.00	0.000000e+00	949.2
1984-01-01	United States of America	9.746000e+02	0.0	0.0	0.00	0.000000e+00	974.6
1992-01-01	United States of America	1.020000e+03	9629090.0	257908.0	26.78	6.540000e+12	1020.0
1996-01-01	United States of America	1.005000e+03	0.0	0.0	0.00	0.000000e+00	1005.0
2002-01-01	United States of America	9.387000e+02	9632030.0	288471.0	29.95	1.100000e+13	938.7

**9b) Create a new column in df\_usa called 'GDP/capita' and populate it with the calculated GDP per capita. Round the results to two decimal points. Display the dataframe head.**

GDP per capita = (GDP / Population)

In [41]:

```
# your code here
df_usa['GDP/capita'] = df_usa['GDP']/df_usa['Population']
df_usa.head()
```

Out[41]:

	Area	Value	Phy Area	Population	Population density	GDP	NRI	GDP/capita
Year								
1962-01-01	United States of America	962909.0	9629090.0	191861.0	19.93	6.050000e+11	0.0	3.153325e+06
1967-01-01	United States of America	962909.0	9629090.0	203713.0	21.16	8.620000e+11	0.0	4.231443e+06
1972-01-01	United States of America	962909.0	9629090.0	213220.0	22.14	1.280000e+12	0.0	6.003189e+06
1977-01-01	United States of America	962909.0	9629090.0	223091.0	23.17	2.090000e+12	0.0	9.368374e+06
1982-01-01	United States of America	962909.0	9629090.0	233954.0	24.30	3.340000e+12	0.0	1.427631e+07

9c) Find the maximum value of the 'NRI' column in the US (using pandas methods). What year does the max value occur? Display the values.

In [42]:

```
# your code here
NRI_max= df_usa['NRI'].max()
df_NRI_max = df_usa[:,df_usa['NRI'] == NRI_max]
print("The max Value of NRI occurs in the following years")
df_NRI_max
```

The max Value of NRI occurs in the following years

Out[42]:

	Area	Value	Phy Area	Population	Population density	GDP	NRI	GDP/c
Year								
1992-01-01	United States of America	9.629090e+05	9629090.0	257908.0	26.78	6.540000e+12	1020.0	2.535788
1992-01-01	United States of America	2.579080e+05	9629090.0	257908.0	26.78	6.540000e+12	1020.0	2.535788
1992-01-01	United States of America	2.678000e+01	9629090.0	257908.0	26.78	6.540000e+12	1020.0	2.535788
1992-01-01	United States of America	6.540000e+12	9629090.0	257908.0	26.78	6.540000e+12	1020.0	2.535788
1992-01-01	United States of America	1.020000e+03	9629090.0	257908.0	26.78	6.540000e+12	1020.0	2.535788

In [ ]: