Colab: https://colab.research.google.com/drive/15Gh2h7l2erNPqfXd7e-Ok8pldEG9kf0N? usp=sharing

```
import pandas as pd
import numpy as np
# GDP per capita V/S Life Expectancy
!wget "https://drive.google.com/uc?export=download&id=1E3bwvYGf1ig32RmcYiWc0IXPN-mD
   ·2022-11-30 15:52:34-- https://drive.google.com/uc?export=download&id=1E3bwvYG
   solving drive.google.com (drive.google.com)... 142.251.2.102, 142.251.2.101, 1
   nnecting to drive.google.com (drive.google.com) | 142.251.2.102 | :443... connecte
   'TP request sent, awaiting response... 303 See Other
   cation: https://doc-0s-68-docs.googleusercontent.com/docs/securesc/ha0ro937gcu
   rning: wildcards not supported in HTTP.
   2022-11-30 15:52:35-- https://doc-0s-68-docs.googleusercontent.com/docs/secur
   solving doc-0s-68-docs.googleusercontent.com (doc-0s-68-docs.googleusercontent
   nnecting to doc-0s-68-docs.googleusercontent.com (doc-0s-68-docs.googleusercon
   'TP request sent, awaiting response... 200 OK
   !ngth: 83785 (82K) [text/csv]
   wing to: 'gapminder.csv'
                                =======>] 81.82K --.-KB/s
                                                                      in 0.001s
 Saving...
                                   'gapminder.csv' saved [83785/83785]
df = pd.read csv("gapminder.csv")
df # structured data, tabular data
```

		country	year	population	continent	life_exp	gdp_cap
	0	Afghanistan	1952	8425333	Asia	28.801	779.445314
type	(df)						
	panda	s.core.fram	e.Data	ıFrame			
	વ	Δfahanietan	1067	11527066	Δeia	3/I U2U	୨ ୧ନ 1071 ୧ ୨
df.s	hape						
	(1704	, 6)					
df["	countr	у"]					
	0	Afghanis	tan				
	1	Afghanis	tan				
	2	Afghanis					
	3	Afghanis					
	4	Afghanis	tan				
	1.600		,				
	1699	Zimba					
	1700	Zimba					
	1701	Zimba					
	1702	Zimba					
	1703	Zimba		1704 d+	o. object		
	name:	country, L	eng ch:	1704, dtyp	e: object		

Saving... ×

df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1704 entries, 0 to 1703
Data columns (total 6 columns):

	,	,	
#	Column	Non-Null Count	Dtype
0	country	1704 non-null	object
1	year	1704 non-null	int64
2	population	1704 non-null	int64
3	continent	1704 non-null	object
4	life_exp	1704 non-null	float64
5	gdp_cap	1704 non-null	float64
dtyp	es: float64(2), int64(2), ob	ject(2)
memo	ry usage: 80	.0+ KB	

df.head(7)

country	year	population	continent	life_exp	gdp_cap	1
Afghanistan	1952	8425333	Asia	28.801	779.445314	
Afghanistan	1957	9240934	Asia	30.332	820.853030	
Afghanistan	1962	10267083	Asia	31.997	853.100710	
Afghanistan	1967	11537966	Asia	34.020	836.197138	
	Afghanistan Afghanistan Afghanistan	Afghanistan 1952 Afghanistan 1957 Afghanistan 1962	Afghanistan 1952 8425333 Afghanistan 1957 9240934 Afghanistan 1962 10267083	Afghanistan 1952 8425333 Asia Afghanistan 1957 9240934 Asia Afghanistan 1962 10267083 Asia	Afghanistan 1957 9240934 Asia 30.332 Afghanistan 1962 10267083 Asia 31.997	Afghanistan 1952 8425333 Asia 28.801 779.445314 Afghanistan 1957 9240934 Asia 30.332 820.853030 Afghanistan 1962 10267083 Asia 31.997 853.100710

df.tail(6)

	country	year	population	continent	life_exp	gdp_cap
1698	Zimbabwe	1982	7636524	Africa	60.363	788.855041
1699	Zimbabwe	1987	9216418	Africa	62.351	706.157306
1700	Zimbabwe	1992	10704340	Africa	60.377	693.420786
1701	Zimbabwe	1997	11404948	Africa	46.809	792.449960
1702	Zimbabwe	2002	11926563	Africa	39.989	672.038623
1703	Zimbabwe	2007	12311143	Africa	43.487	469.709298

row orineted approach - lists of lists

```
pd.DataFrame([['Afghanistan',1952, 8425333, 'Asia', 28.801, 779.445314 ],

Saving...

Saving...

year','population','continent','life_exp','gdp_c
```

	country	year	population	continent	life_exp	gdp_cap	1
0	Afghanistan	1952	8425333	Asia	28.801	779.445314	
1	Afghanistan	1957	9240934	Asia	30.332	820.853030	
2	Afghanistan	1962	102267083	Asia	31.997	853.100710	

	country	year	population	continent	life_exp	gdp_cap	1
0	Afghanistan	1952	8425333	Asia	28.801	779.445314	

column oriented approach

	country	year	population	continent	life_exp	gdp_cap	10-
0	Afghanistan	1952	842533	Asia	28.801	779.445314	
1	Afghanistan	1957	9240934	Asia	30.332	820.853030	

	country	year	population	continent	life_exp	1
0	Afghanistan	1952	842533	779.445314	28.801	
1	Afghanistan	1957	9240934	820.853030	30.332	

df.columns

Index(['country', 'year', 'population', 'continent', 'life_exp', 'gdp_cap'],
dtype='object')

df.keys() # df works like a "specialised" dictionary

Index(['country', 'year', 'population', 'continent', 'life_exp', 'gdp_cap'],
dtype='object')

Saving...

X ract subset of columns

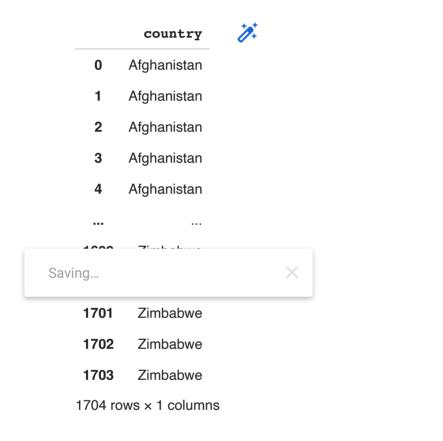
	country	life_exp
0	Afghanistan	28.801
1	Afghanistan	30.332
2	Afghanistan	31.997
3	Afghanistan	34.020
4	Afghanistan	36.088
1699	Zimbabwe	62.351
1700	Zimbabwe	60.377
1701	Zimbabwe	46.809
1702	Zimbabwe	39.989
1703	Zimbabwe	43.487
1704	0	_

1704 rows × 2 columns

df["country"] # series

```
Afghanistan
1
        Afghanistan
2
        Afghanistan
3
        Afghanistan
4
        Afghanistan
           . . .
1699
           Zimbabwe
1700
           Zimbabwe
1701
           Zimbabwe
1702
           Zimbabwe
1703
           Zimbabwe
Name: country, Length: 1704, dtype: object
```

df[["country"]] # returns as dataframe because of double brackers



df["country"].unique() #np.unique(df["country"])

'Korea, Rep.', 'Kuwait', 'Lebanon', 'Lesotho', 'Liberia', 'Libya',

```
'Madagascar', 'Malawi', 'Malaysia', 'Mali', 'Mauritania',
'Mauritius', 'Mexico', 'Mongolia', 'Montenegro', 'Morocco',
'Mozambique', 'Myanmar', 'Namibia', 'Nepal', 'Netherlands',
'New Zealand', 'Nicaragua', 'Niger', 'Nigeria', 'Norway', 'Oman',
'Pakistan', 'Panama', 'Paraguay', 'Peru', 'Philippines', 'Poland',
'Portugal', 'Puerto Rico', 'Reunion', 'Romania', 'Rwanda',
'Sao Tome and Principe', 'Saudi Arabia', 'Senegal', 'Serbia',
'Sierra Leone', 'Singapore', 'Slovak Republic', 'Slovenia',
'Somalia', 'South Africa', 'Spain', 'Sri Lanka', 'Sudan',
'Swaziland', 'Sweden', 'Switzerland', 'Syria', 'Taiwan',
'Tanzania', 'Thailand', 'Togo', 'Trinidad and Tobago', 'Tunisia',
'Turkey', 'Uganda', 'United Kingdom', 'United States', 'Uruguay',
'Venezuela', 'Vietnam', 'West Bank and Gaza', 'Yemen, Rep.',
'Zambia', 'Zimbabwe'], dtype=object)
```

df["country"].nunique()

142

df["country"].value counts()

Afghanistan 12
Pakistan 12
New Zealand 12
Nicaragua 12
Niger 12
Eritrea 12

Saving... X

Name: country, Length: 142, dtype: int64

df.rename({"population":"Population", "country":"Country"}, axis=1, inplace=True)

df

	country	year	population	continent	life_exp	gdp_cap	10-
0	Afghanistan	1952	8425333	Asia	28.801	779.445314	
1	Afghanistan	1957	9240934	Asia	30.332	820.853030	

df

	Country	year	Population	continent	life_exp	gdp_cap
0	Afghanistan	1952	8425333	Asia	28.801	779.445314
1	Afghanistan	1957	9240934	Asia	30.332	820.853030
2	Afghanistan	1962	10267083	Asia	31.997	853.100710
3	Afghanistan	1967	11537966	Asia	34.020	836.197138
4	Afghanistan	1972	13079460	Asia	36.088	739.981106
	•••					
1699	Zimbabwe	1987	9216418	Africa	62.351	706.157306
1700	Zimbabwe	1992	10704340	Africa	60.377	693.420786
1701	Zimbabwe	1997	11404948	Africa	46.809	792.449960
1702	Zimbabwe	2002	11926563	Africa	39.989	672.038623
1703	Zimbabwe	2007	12311143	Africa	43.487	469.709298
ing			×			

df.rename({"continent":"Continent"}) # default axis=0

	Country	year	Population	continent	life_exp	gdp_cap	
0	Afghanistan	1952	8425333	Asia	28.801	779.445314	
1	Afghanistan	1957	9240934	Asia	30.332	820.853030	
2	Afghanistan	1962	10267083	Asia	31.997	853.100710	
3	Afghanistan	1967	11537966	Asia	34.020	836.197138	
4	Afghanistan	1972	13079460	Asia	36.088	739.981106	
1699	Zimbabwe	1987	9216418	Africa	62.351	706.157306	
1700	Zimbabwe	1992	10704340	Africa	60.377	693.420786	
1701	Zimbabwe	1997	11404948	Africa	46.809	792.449960	
1702	Zimbabwe	2002	11926563	Africa	39.989	672.038623	
1703	Zimbabwe	2007	12311143	Africa	43.487	469.709298	
1704							

1704 rows × 6 columns

```
df["Country"] # dict like style
```

```
0
        Afghanistan
1
        Afghanistan
2
        Afghanistan
3
        Afghanistan
        Afghanistan
           . . .
1699
           Zimbabwe
           Zimbabwe
1700
1701
           Zimbabwe
           Zimbabwe
1702
           Zimbabwe
1703
Name: Country, Length: 1704, dtype: object
```

df.Country # attribute, DONT USE THIS

0	Afghanistan			
1	Afghanistan			
2	Afghanistan			
3	Afghanistan			
4	Afghanistan			
	• • •			
1699	Zimbabwe			
1700	Zimbabwe			
1701	Zimbabwe			
1702	Zimbabwe			
1703	7 i mhahwo			
Saving		X	dtype:	object

HOMEWORK - WHY I SHOULD NOT USE THIS STYLE

df.drop("continent", axis=1) # for permanent changes use inplace=True

	Country	year	Population	life_exp	gdp_cap
-					
df["year+	-7"] = df["ye	ear"] -	+ 7		
	J				
df["gdp"]	= df["gdp_c	cap"]	* df["Popula	tion"]	
3	Afghanistan	1967	11537966	34.020	836.197138
df					

	Country	year	Population	continent	life_exp	gdp_cap	gdp
0	Afghanistan	1952	8425333	Asia	28.801	779.445314	6.567086e+09
1	Afghanistan	1957	9240934	Asia	30.332	820.853030	7.585449e+09
2	Afghanistan	1962	10267083	Asia	31.997	853.100710	8.758856e+09
3	Afghanistan	1967	11537966	Asia	34.020	836.197138	9.648014e+09
4	Afghanistan	1972	13079460	Asia	36.088	739.981106	9.678553e+09
1699	Zimbabwe	1987	9216418	Africa	62.351	706.157306	6.508241e+09
1700	Zimbabwe	1992	10704340	Africa	60.377	693.420786	7.422612e+09
1701	Zimbabwe	1997	11404948	Africa	46.809	792.449960	9.037851e+09
1700	7:mbaba	0000	11026563	Africa	39.989	672.038623	8.015111e+09
Saving			× 11143	Africa	43.487	469.709298	5.782658e+09

1704 rows × 7 columns

df["Own"] = [i for i in range(len(df))]

df

	Country	year	Population	continent	life_exp	gdp_cap	gdp
0	Afghanistan	1952	8425333	Asia	28.801	779.445314	6.567086e+09
1	Afghanistan	1957	9240934	Asia	30.332	820.853030	7.585449e+09

df.insert() that helps in adding a column at a particular location

3 AIYIIAIIISIAII 1307 11037300 ASIA 34.020 030.137130 3.0400146403

Working with Rows

```
ser = df["Country"]
ser
```

0	Afghanistan
1	Afghanistan
2	Afghanistan
3	Afghanistan
4	Afghanistan
	• • •
1699	Zimbabwe
1700	Zimbabwe
1701	7 i mb a brea

1701 Zimbabwe 1702 Zimbabwe 1703 Zimbabwe

Name: Country, Length: 1704, dtype: object



ser[5:14]

- 5 Afghanistan
- 6 Afghanistan
- 7 Afghanistan
- 8 Afghanistan
- 9 Afghanistan
- 10 Afghanistan
- 11 Afghanistan
- 12 Albania
- 13 Albania

Name: Country, dtype: object

df[0] # indexing a row like doesnt happan in dataframe # df["country"] #looks for this index along axis=1

```
_____
```

3360 try:

-> 3361 return self._engine.get_loc(casted_key)
3362 except KeyError as err:

4 frames -

pandas/_libs/hashtable_class_helper.pxi in pandas._libs.hashtable.PyObjectHashTable.get_item()

pandas/_libs/hashtable_class_helper.pxi in pandas._libs.hashtable.PyObjectHashTable.get_item()

KeyError: 0

The above exception was the direct cause of the following exception:

KeyError

Traceback (most recent call last)

/usr/local/lib/python3.7/dist-packages/pandas/core/indexes/base.py in
get_loc(self, key, method, tolerance)

return self. engine.get loc(casted kev)

df[5:14]

		Country	year	Population	continent	life_exp	gdp_cap	gdp	0
	5	Afghanistan	1977	14880372	Asia	38.438	786.113360	1.169766e+10	
	0			316	Asia	39.854	978.011439	1.259856e+10	
L	Saving			× 957	Asia	40.822	852.395945	1.182099e+10	
	8	Afghanistan	1992	16317921	Asia	41.674	649.341395	1.059590e+10	
	9	Afghanistan	1997	22227415	Asia	41.763	635.341351	1.412200e+10	
	10	Afghanistan	2002	25268405	Asia	42.129	726.734055	1.836341e+10	
	11	Afghanistan	2007	31889923	Asia	43.828	974.580338	3.107929e+10	
	12	Albania	1952	1282697	Europe	55.230	1601.056136	2.053670e+09	
	13	Albania	1957	1476505	Europe	59.280	1942.284244	2.867792e+09	

===> Indexing doesnt work to index rows because of similar syntax for columns
===> Slicing works for slicing the rows

df.index.values

```
array([ 0, 1, 2, ..., 1701, 1702, 1703])
```

df.columns

df.index = range(1, df.shape[0]+1)

df

	Country	year	Population	continent	life_exp	gdp_cap	gdp
1	Afghanistan	1952	8425333	Asia	28.801	779.445314	6.567086e+09
2	Afghanistan	1957	9240934	Asia	30.332	820.853030	7.585449e+09
3	Afghanistan	1962	10267083	Asia	31.997	853.100710	8.758856e+09
4	Afghanistan	1967	11537966	Asia	34.020	836.197138	9.648014e+09
5	Afghanistan	1972	13079460	Asia	36.088	739.981106	9.678553e+09
1700	Zimbabwe	1987	9216418	Africa	62.351	706.157306	6.508241e+09
1701	Zimbabwe	1992	10704340	Africa	60.377	693.420786	7.422612e+09
1702	Zimbabwe	1997	11404948	Africa	46.809	792.449960	9.037851e+09
1703	Zimbabwe	2002	11926563	Africa	39.989	672.038623	8.015111e+09
1704	Zimbabwe	2007	12311143	Africa	43.487	469.709298	5.782658e+09

1704 rows × 8 columns



2

df.index = np.arange(1, df.shape[0]+1, dtype='float')
df

		Country	year	Population	continent	lite_exp	gap_cap	gap
	1.0	Afghanistan	1952	8425333	Asia	28.801	779.445314	6.567086e+09
sampl	le = df	.head()						
	٠.٠	,		.020,000	, 1014	JJ.	0000070	0.7000000100

sample

	Country	year	Population	continent	life_exp	gdp_cap	gdp	70
1.0	Afghanistan	1952	8425333	Asia	28.801	779.445314	6.567086e+09	
2.0	Afghanistan	1957	9240934	Asia	30.332	820.853030	7.585449e+09	
3.0	Afghanistan	1962	10267083	Asia	31.997	853.100710	8.758856e+09	
4.0	Afghanistan	1967	11537966	Asia	34.020	836.197138	9.648014e+09	
5.0	Afghanistan	1972	13079460	Asia	36.088	739.981106	9.678553e+09	

sample.index = ["a", "b", "c", "d", "e"]

sample

Own	gdp	gdp_cap	life_exp	continent	Population	year	Country		
С	6.567086e+09	779.445314	28.801	Asia	33			Saving	
1	7.585449e+09	820.853030	30.332	Asia	34		_	Saviriy.	
2	8.758856e+09	853.100710	31.997	Asia	10267083	1962	Afghanistan	С	
3	9.648014e+09	836.197138	34.020	Asia	11537966	1967	Afghanistan	d	
4	9.678553e+09	739.981106	36.088	Asia	13079460	1972	Afghanistan	е	

sample.index = ["a", "b", "c", "d", "d"]

sample

	Country	year	Population	continent	life_exp	gdp_cap	gdp	Own
а	Afghanistan	1952	8425333	Asia	28.801	779.445314	6.567086e+09	С
b	Afghanistan	1957	9240934	Asia	30.332	820.853030	7.585449e+09	1
С	Afghanistan	1962	10267083	Asia	31.997	853.100710	8.758856e+09	2
d	Afghanistan	1967	11537966	Asia	34.020	836.197138	9.648014e+09	3
d	Afghanistan	1972	13079460	Asia	36.088	739.981106	9.678553e+09	4

df.columns

df

		Country	year	Population	continent	life_exp	gdp_cap	gdp
1.	0	Afghanistan	1952	8425333	Asia	28.801	779.445314	6.567086e+09
2.	0	Afghanistan	1957	9240934	Asia	30.332	820.853030	7.585449e+09
3.	0	Afghanistan	1962	10267083	Asia	31.997	853.100710	8.758856e+09
4.	0	Afghanistan	1967	11537966	Asia	34.020	836.197138	9.648014e+09
5.	0	Afghanistan	1972	13079460	Asia	36.088	739.981106	9.678553e+09
170	0.0	Zimbabwe	1987	9216418	Africa	62.351	706.157306	6.508241e+09
170	1.0	Zimbabwe	1992	10704340	Africa	60.377	693.420786	7.422612e+09
170	2.0	Zimbabwe	1997	11404948	Africa	46.809	792.449960	9.037851e+09
Saving				926563	Africa	39.989	672.038623	8.015111e+09
Saving				2311143	Africa	43.487	469.709298	5.782658e+09

1704 rows × 8 columns

duplication of explciit indexes is ok

df["Country"] # you can basically group the data, classification using same explcit

	Country	Country	1	
1.0	Afghanistan	0		
2.0	Afghanistan	1		
2.0	Afahaniatan	0		

df.reset_index() # drop=True would have dropped the colum

		index	Country	year	Population	continent	life_exp	gdp_cap	
	0	1.0	Afghanistan	1952	8425333	Asia	28.801	779.445314	6.56708
	1	2.0	Afghanistan	1957	9240934	Asia	30.332	820.853030	7.58544
	2	3.0	Afghanistan	1962	10267083	Asia	31.997	853.100710	8.75885
	3	4.0	Afghanistan	1967	11537966	Asia	34.020	836.197138	9.64801
	4	5.0	Afghanistan	1972	13079460	Asia	36.088	739.981106	9.67855
	1699	1700.0	Zimbabwe	1987	9216418	Africa	62.351	706.157306	6.50824
	1700	1701.0	Zimbabwe	1992	10704340	Africa	60.377	693.420786	7.42261
	1701	1702.0	Zimbabwe	1997	11404948	Africa	46.809	792.449960	9.03785
	1702	1703.0	Zimbabwe	2002	11926563	Africa	39.989	672.038623	8.01511
	1703	1704.0	Zimbabwe	2007	12311143	Africa	43.487	469.709298	5.78265
Savi	ng			×					

Learner's personal doubt

Double-click (or enter) to edit

```
x = pd.DataFrame({'veclocity':[100, -11, -16, 13, 14]})

def find_closest(x):
    return np.argmin(np.abs(np.array([3, 30, -5, -15]) - x)))

x["veclocity"].apply(find_closest)

    0    1
    1    3
    2    3
    3    0
    4    0
    Name: veclocity, dtype: int64
```

Colab paid products - Cancel contracts here

completed at 23:26

✓ 0s

Saving... ×

×