

# Pandas

## Introduction:

- Name pandas is derived from the word **Panel Data** – an econometrics term for multidimensional data
- It is an Open-source Python library that provides high performance data manipulation and analysis using its powerful data structures
- It has its own s high-performance, easy-to-use data structures and analysis tools for the Python programming language
- It deals with **2D data** called as **Data Frame**
- Usually, it requires following import files:  
**import pandas as pd**  
**import numpy as np**  
**import os**
- We need to set path of directory where input files are kept using:  
**os.chdir('D:/Pandas')**

Sr No	Method	Example and Output	Explanation																																																																																												
1.	Dataframe=pd.read_csv('file_name.csv')	<div><div>a) Importing Data</div><div><pre>cars_data = pd.read_csv('Toyota.csv')</pre></div><div><table><tr><th>Index</th><th>Unnamed: 0</th><th>Price</th><th>Age</th><th>KM</th><th>FuelType</th><th>HP</th><th>MetColor</th><th>Automatic</th><th>CC</th><th>Doors</th><th>Weight</th></tr><tr><td>0</td><td>0</td><td>13500</td><td>23</td><td>46986</td><td>Diesel</td><td>90</td><td>1</td><td>0</td><td>2000</td><td>three</td><td>1165</td></tr><tr><td>1</td><td>1</td><td>13750</td><td>23</td><td>72937</td><td>Diesel</td><td>90</td><td>1</td><td>0</td><td>2000</td><td>3</td><td>1165</td></tr><tr><td>2</td><td>2</td><td>13950</td><td>24</td><td>41711</td><td>Diesel</td><td>90</td><td>nan</td><td>0</td><td>2000</td><td>3</td><td>1165</td></tr></table></div><div><div>b) Importing Data by setting first column as index column</div><div><pre>cars_data = pd.read_csv('Toyota.csv',index_col=0 )</pre></div><div><table><tr><th>Index</th><th>Price</th><th>Age</th><th>KM</th><th>FuelType</th><th>HP</th><th>MetColor</th><th>Automatic</th><th>CC</th><th>Doors</th><th>Weight</th></tr><tr><td>0</td><td>13500</td><td>23</td><td>46986</td><td>Diesel</td><td>90</td><td>1</td><td>0</td><td>2000</td><td>three</td><td>1165</td></tr><tr><td>1</td><td>13750</td><td>23</td><td>72937</td><td>Diesel</td><td>90</td><td>1</td><td>0</td><td>2000</td><td>3</td><td>1165</td></tr><tr><td>2</td><td>13950</td><td>24</td><td>41711</td><td>Diesel</td><td>90</td><td>nan</td><td>0</td><td>2000</td><td>3</td><td>1165</td></tr></table></div></div></div>	Index	Unnamed: 0	Price	Age	KM	FuelType	HP	MetColor	Automatic	CC	Doors	Weight	0	0	13500	23	46986	Diesel	90	1	0	2000	three	1165	1	1	13750	23	72937	Diesel	90	1	0	2000	3	1165	2	2	13950	24	41711	Diesel	90	nan	0	2000	3	1165	Index	Price	Age	KM	FuelType	HP	MetColor	Automatic	CC	Doors	Weight	0	13500	23	46986	Diesel	90	1	0	2000	three	1165	1	13750	23	72937	Diesel	90	1	0	2000	3	1165	2	13950	24	41711	Diesel	90	nan	0	2000	3	1165	Imports the data of .csv file into dataframe
Index	Unnamed: 0	Price	Age	KM	FuelType	HP	MetColor	Automatic	CC	Doors	Weight																																																																																				
0	0	13500	23	46986	Diesel	90	1	0	2000	three	1165																																																																																				
1	1	13750	23	72937	Diesel	90	1	0	2000	3	1165																																																																																				
2	2	13950	24	41711	Diesel	90	nan	0	2000	3	1165																																																																																				
Index	Price	Age	KM	FuelType	HP	MetColor	Automatic	CC	Doors	Weight																																																																																					
0	13500	23	46986	Diesel	90	1	0	2000	three	1165																																																																																					
1	13750	23	72937	Diesel	90	1	0	2000	3	1165																																																																																					
2	13950	24	41711	Diesel	90	nan	0	2000	3	1165																																																																																					

2.	DataFrame.index	<pre>In [14]: import os import pandas as pd import numpy as np os.chdir("Z:/Pandas") car=pd.read_csv('Toyota.csv',index_col=0)</pre> <hr/> <pre>In [13]: car.index</pre> <hr/> <pre>Out[13]: Int64Index([ 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, ... 1426, 1427, 1428, 1429, 1430, 1431, 1432, 1433, 1434, 1435], dtype='int64', length=1436)</pre>	To get the index (row labels) of the dataframe
3.	Dataframe.columns	<pre>In [12]: car.columns</pre> <hr/> <pre>Out[12]: Index(['Price', 'Age', 'KM', 'FuelType', 'HP', 'MetColor', 'Automatic', 'CC', 'Doors', 'Weight'], dtype='object')</pre>	To get all the column labels
4.	Dataframe.size	<pre>In [15]: car.size</pre> <hr/> <pre>Out[15]: 14360</pre>	To get the total number of elements from the data frame
5.	Dataframe.shape	<pre>In [16]: car.shape</pre> <hr/> <pre>Out[16]: (1436, 10)</pre>	To get the shape of the data frame
6.	Dataframe.memory_usage	<pre>car.memory_usage()</pre> <hr/> <pre>Out[17]: Index      11488 Price      11488 Age        11488 KM         11488 FuelType   11488 HP         11488 MetColor   11488 Automatic   11488 CC         11488 Doors      11488 Weight     11488 dtype: int64</pre>	To get Memory occupied by each column in bytes
7.	Dataframe.ndim	<pre>In [18]: car.ndim</pre> <hr/> <pre>Out[18]: 2</pre>	To get the Dimension of data frame

8.	Dataframe.head([n])	<div><pre>car.head(5)</pre></div> <div>Out[19]:</div> <table><thead><tr><th></th><th>Price</th><th>Age</th><th>KM</th><th>FuelType</th><th>HP</th><th>MetColor</th><th>Automatic</th><th>CC</th><th>Doors</th><th>Weight</th></tr></thead><tbody><tr><td>0</td><td>13500</td><td>23.0</td><td>46986</td><td>Diesel</td><td>90</td><td>1.0</td><td>0</td><td>2000</td><td>three</td><td>1165</td></tr><tr><td>1</td><td>13750</td><td>23.0</td><td>72937</td><td>Diesel</td><td>90</td><td>1.0</td><td>0</td><td>2000</td><td>3</td><td>1165</td></tr><tr><td>2</td><td>13950</td><td>24.0</td><td>41711</td><td>Diesel</td><td>90</td><td>NaN</td><td>0</td><td>2000</td><td>3</td><td>1165</td></tr><tr><td>3</td><td>14950</td><td>26.0</td><td>48000</td><td>Diesel</td><td>90</td><td>0.0</td><td>0</td><td>2000</td><td>3</td><td>1165</td></tr><tr><td>4</td><td>13750</td><td>30.0</td><td>38500</td><td>Diesel</td><td>90</td><td>0.0</td><td>0</td><td>2000</td><td>3</td><td>1170</td></tr></tbody></table>		Price	Age	KM	FuelType	HP	MetColor	Automatic	CC	Doors	Weight	0	13500	23.0	46986	Diesel	90	1.0	0	2000	three	1165	1	13750	23.0	72937	Diesel	90	1.0	0	2000	3	1165	2	13950	24.0	41711	Diesel	90	NaN	0	2000	3	1165	3	14950	26.0	48000	Diesel	90	0.0	0	2000	3	1165	4	13750	30.0	38500	Diesel	90	0.0	0	2000	3	1170	The function head returns the first n rows from the dataframe
	Price	Age	KM	FuelType	HP	MetColor	Automatic	CC	Doors	Weight																																																											
0	13500	23.0	46986	Diesel	90	1.0	0	2000	three	1165																																																											
1	13750	23.0	72937	Diesel	90	1.0	0	2000	3	1165																																																											
2	13950	24.0	41711	Diesel	90	NaN	0	2000	3	1165																																																											
3	14950	26.0	48000	Diesel	90	0.0	0	2000	3	1165																																																											
4	13750	30.0	38500	Diesel	90	0.0	0	2000	3	1170																																																											
9.	Dataframe.tail(n)	<div><pre>car.tail(5)</pre></div> <div>Out[20]:</div> <table><thead><tr><th></th><th>Price</th><th>Age</th><th>KM</th><th>FuelType</th><th>HP</th><th>MetColor</th><th>Automatic</th><th>CC</th><th>Doors</th><th>Weight</th></tr></thead><tbody><tr><td>1431</td><td>7500</td><td>NaN</td><td>20544</td><td>Petrol</td><td>86</td><td>1.0</td><td>0</td><td>1300</td><td>3</td><td>1025</td></tr><tr><td>1432</td><td>10845</td><td>72.0</td><td>??</td><td>Petrol</td><td>86</td><td>0.0</td><td>0</td><td>1300</td><td>3</td><td>1015</td></tr><tr><td>1433</td><td>8500</td><td>NaN</td><td>17016</td><td>Petrol</td><td>86</td><td>0.0</td><td>0</td><td>1300</td><td>3</td><td>1015</td></tr><tr><td>1434</td><td>7250</td><td>70.0</td><td>??</td><td>NaN</td><td>86</td><td>1.0</td><td>0</td><td>1300</td><td>3</td><td>1015</td></tr><tr><td>1435</td><td>6950</td><td>76.0</td><td>1</td><td>Petrol</td><td>110</td><td>0.0</td><td>0</td><td>1600</td><td>5</td><td>1114</td></tr></tbody></table>		Price	Age	KM	FuelType	HP	MetColor	Automatic	CC	Doors	Weight	1431	7500	NaN	20544	Petrol	86	1.0	0	1300	3	1025	1432	10845	72.0	??	Petrol	86	0.0	0	1300	3	1015	1433	8500	NaN	17016	Petrol	86	0.0	0	1300	3	1015	1434	7250	70.0	??	NaN	86	1.0	0	1300	3	1015	1435	6950	76.0	1	Petrol	110	0.0	0	1600	5	1114	Returns Last n Rows
	Price	Age	KM	FuelType	HP	MetColor	Automatic	CC	Doors	Weight																																																											
1431	7500	NaN	20544	Petrol	86	1.0	0	1300	3	1025																																																											
1432	10845	72.0	??	Petrol	86	0.0	0	1300	3	1015																																																											
1433	8500	NaN	17016	Petrol	86	0.0	0	1300	3	1015																																																											
1434	7250	70.0	??	NaN	86	1.0	0	1300	3	1015																																																											
1435	6950	76.0	1	Petrol	110	0.0	0	1600	5	1114																																																											
10	Dataframe.at[index, Label]	<div><pre>car.at[4, 'KM']</pre></div> <div>Out[22]: '38500'</div>	Label based Scalar Lookup																																																																		
11.	Dataframe.iat[index, index]	<div><pre>car.iat[4,9]</pre></div> <div>Out[24]: 1170</div>	Index based Lookup																																																																		
12.	Dataframe.loc[Row : Row, Col : Col ]	<div>a)</div> <div><pre>car.loc[:, 'FuelType']</pre></div> <div>Out[26]: 0 Diesel 1 Diesel 2 Diesel 3 Diesel 4 Diesel ... 1431 Petrol 1432 Petrol 1433 Petrol 1434 NaN 1435 Petrol Name: FuelType, Length: 1436, dtype: object</div>	Accessing Group of rows and Columns by Labels																																																																		

b)

```
car.loc[0:10, 'FuelType']
```

```
Out[27]: 0    Diesel
         1    Diesel
         2    Diesel
         3    Diesel
         4    Diesel
         5    Diesel
         6    Diesel
         7     NaN
         8   Petrol
         9    Diesel
        10   Petrol
        Name: FuelType, dtype: object
```

c)

```
In [28]: #Accessing all the elements of row 1
```

```
car.loc[1,:]
```

```
Out[28]: Price      13750
         Age         23
         KM       72937
         FuelType   Diesel
         HP         90
         MetColor    1
         Automatic    0
         CC       2000
         Doors        3
         Weight    1165
        Name: 1, dtype: object
```

d)

```
In [34]: #row 0 column Price to HP
```

```
car.loc[0, 'Price': 'HP']
```

```
Out[34]: Price      13500
         Age         23
         KM      46986
         FuelType   Diesel
         HP         90
        Name: 0, dtype: object
```

## Pandas Data Types and Data Cleansing:

- The way information gets stored in a dataframe or a python object affects the analysis and outputs of calculations
- There are two main types of data
  - numeric

Python data type	Pandas data type	Description
int	int64	Numeric characters
float	float64	Numeric characters with decimals

- character types
  - category & object

category	object
<ul style="list-style-type: none"><li>• A string variable consisting of only a few different values. Converting such a string variable to a categorical variable will save some memory</li><li>• A categorical variable takes on a limited, fixed number of possible values</li></ul>	<ul style="list-style-type: none"><li>• The column will be assigned as object data type when it has mixed types (numbers and strings). If a column contains 'nan'(blank cells), pandas will default to object datatype.</li><li>• For strings, the length is not fixed</li></ul>

- Numeric data types includes integers and floats
  - For example: integer – 10, float – 10.53
- Strings are known as objects in pandas which can store values that contain numbers and / or characters
  - For example: 'category1'

### Learning Outcome of Point 13-17:

- Data types ○ Numeric ○ Character
- Checked data types of each column
- Count of unique data types
- Selected data based on data types
- Concise summary of dataframe
- Checked format of each column
- Got unique elements of each column

Sr No	Method	Example and Output	Explanation																																																																																				
13.	DataFrame.dtypes	<pre>car.dtypes</pre> <pre>Out[4]: Price      int64 Age        float64 KM         object FuelType   object HP         object MetColor   float64 Automatic  int64 CC         int64 Doors      object Weight     int64 dtype: object</pre>	dtypes returns a series with the data type of each column																																																																																				
14.	DataFrame.get_dtype_counts()	<pre>In [15]: #3) count of unique data types</pre> <pre>car.dtypes.value_counts()</pre> <pre>Out[15]: object      4 int64      4 float64     2 dtype: int64</pre>	It returns counts of unique data types in the dataframe																																																																																				
15.	DataFrame.select_dtypes(include=None, exclude=None)	<pre>In [17]: #4) selecting data based on data type</pre> <pre>car.select_dtypes(exclude=[object])</pre> <pre>Out[17]:</pre> <table><thead><tr><th></th><th>Price</th><th>Age</th><th>MetColor</th><th>Automatic</th><th>CC</th><th>Weight</th></tr></thead><tbody><tr><td>0</td><td>13500</td><td>23.0</td><td>1.0</td><td>0</td><td>2000</td><td>1165</td></tr><tr><td>1</td><td>13750</td><td>23.0</td><td>1.0</td><td>0</td><td>2000</td><td>1165</td></tr><tr><td>2</td><td>13950</td><td>24.0</td><td>NaN</td><td>0</td><td>2000</td><td>1165</td></tr><tr><td>3</td><td>14950</td><td>26.0</td><td>0.0</td><td>0</td><td>2000</td><td>1165</td></tr><tr><td>4</td><td>13750</td><td>30.0</td><td>0.0</td><td>0</td><td>2000</td><td>1170</td></tr><tr><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td></tr><tr><td>1431</td><td>7500</td><td>NaN</td><td>1.0</td><td>0</td><td>1300</td><td>1025</td></tr><tr><td>1432</td><td>10845</td><td>72.0</td><td>0.0</td><td>0</td><td>1300</td><td>1015</td></tr><tr><td>1433</td><td>8500</td><td>NaN</td><td>0.0</td><td>0</td><td>1300</td><td>1015</td></tr><tr><td>1434</td><td>7250</td><td>70.0</td><td>1.0</td><td>0</td><td>1300</td><td>1015</td></tr><tr><td>1435</td><td>6950</td><td>76.0</td><td>0.0</td><td>0</td><td>1600</td><td>1114</td></tr></tbody></table> <p>1436 rows × 6 columns</p>		Price	Age	MetColor	Automatic	CC	Weight	0	13500	23.0	1.0	0	2000	1165	1	13750	23.0	1.0	0	2000	1165	2	13950	24.0	NaN	0	2000	1165	3	14950	26.0	0.0	0	2000	1165	4	13750	30.0	0.0	0	2000	1170	...	...	...	...	...	...	...	1431	7500	NaN	1.0	0	1300	1025	1432	10845	72.0	0.0	0	1300	1015	1433	8500	NaN	0.0	0	1300	1015	1434	7250	70.0	1.0	0	1300	1015	1435	6950	76.0	0.0	0	1600	1114	It returns a subset of the columns from dataframe based on the column dtypes
	Price	Age	MetColor	Automatic	CC	Weight																																																																																	
0	13500	23.0	1.0	0	2000	1165																																																																																	
1	13750	23.0	1.0	0	2000	1165																																																																																	
2	13950	24.0	NaN	0	2000	1165																																																																																	
3	14950	26.0	0.0	0	2000	1165																																																																																	
4	13750	30.0	0.0	0	2000	1170																																																																																	
...	...	...	...	...	...	...																																																																																	
1431	7500	NaN	1.0	0	1300	1025																																																																																	
1432	10845	72.0	0.0	0	1300	1015																																																																																	
1433	8500	NaN	0.0	0	1300	1015																																																																																	
1434	7250	70.0	1.0	0	1300	1015																																																																																	
1435	6950	76.0	0.0	0	1600	1114																																																																																	

16.	info()	<pre>In [18]: #5) Concise Summary of dataframe and 6) Checking Format of Each Column  car.info()  &lt;class 'pandas.core.frame.DataFrame'&gt; Int64Index: 1436 entries, 0 to 1435 Data columns (total 10 columns): #   Column      Non-Null Count  Dtype ---  --- 0   Price       1436 non-null   int64 1   Age         1336 non-null   float64 2   KM          1436 non-null   object 3   FuelType    1336 non-null   object 4   HP          1436 non-null   object 5   MetColor    1286 non-null   float64 6   Automatic   1436 non-null   int64 7   CC          1436 non-null   int64 8   Doors       1436 non-null   object 9   Weight      1436 non-null   int64 dtypes: float64(2), int64(4), object(4) memory usage: 123.4+ KB</pre>	<p>returns a concise summary of a dataframe</p> <ul style="list-style-type: none"><li>• data type of index</li><li>• data type of columns</li><li>• count of non-null values</li><li>• memory usage</li></ul>
17.	numpy.unique(array)	<pre>In [19]: #7) Getting Unique Elements of each Column  print(np.unique(car['KM']))  ['1' '10000' '100123' ... '99865' '99971' '??']  In [20]: print(np.unique(car['HP']))  ['107' '110' '116' '192' '69' '71' '72' '73' '86' '90' '97' '98' '????']  In [21]: print(np.unique(car['MetColor']))  [ 0.  1. nan nan nan nan nan nan nan nan nan nan nan nan nan nan nan nan  nan nan nan nan nan nan nan nan nan nan nan nan nan nan nan nan nan  nan nan nan nan nan nan nan nan nan nan nan nan nan nan nan nan nan  nan nan nan nan nan nan nan nan nan nan nan nan nan nan nan nan nan  nan nan nan nan nan nan nan nan nan nan nan nan nan nan nan nan nan  nan nan nan nan nan nan nan nan nan nan nan nan nan nan nan nan nan  nan nan nan nan nan nan nan nan nan nan nan nan nan nan nan nan nan  nan nan nan nan nan nan nan nan nan]  In [22]: print(np.unique(car['Automatic']))  [0 1]  In [23]: print(np.unique(car['Doors']))  ['2' '3' '4' '5' 'five' 'four' 'three']</pre>	<p><i>unique()</i> is used to find the unique elements of a column</p> <p>‘HP’ has special character to it – Hence, it has been read as object instead of int64</p> <p>‘Doors’ has been read as object instead of int64 because of values ‘five’ ‘four’ ‘three’ which are strings</p>

## Missing Values:

We need to know how missing values are represented in the dataset in order to make reasonable decisions

The missing values exist in the form of 'nan' '??', '????'

Python, by default replace blank values with 'nan'

Now, importing the data considering other forms of missing values in a dataframe

□ Importing data □ Concise summary of dataframe □ Converting variable's data types □ Category vs Object data type □ Cleaning column 'Doors' □ Getting count of missing values

Sr No	Method	Example and Output	Explanation
18	Dataframe=pd.read_csv('file_name.csv', na_values=['??','????'])	<pre>car=pd.read_csv('Toyota.csv', index_col=0, na_values=["??", "????"]) car.info()  &lt;class 'pandas.core.frame.DataFrame'&gt; Int64Index: 1436 entries, 0 to 1435 Data columns (total 10 columns): #   Column      Non-Null Count  Dtype ---  - 0   Price       1436 non-null  int64 1   Age         1336 non-null  float64 2   KM          1421 non-null  float64 3   FuelType    1336 non-null  object 4   HP          1430 non-null  float64 5   MetColor    1286 non-null  float64 6   Automatic   1436 non-null  int64 7   CC          1436 non-null  int64 8   Doors       1436 non-null  object 9   Weight      1436 non-null  int64 dtypes: float64(4), int64(4), object(2) memory usage: 123.4+ KB</pre>	Imports the data of .csv file into dataframe by considering other forms of data



19	<div>Dataframe.astype( )</div> <div>Dataframe.nbytes( )</div>	<pre>car['MetColor']=car['MetColor'].astype('object') car['Automatic']=car['Automatic'].astype('object')  car.info()  &lt;class 'pandas.core.frame.DataFrame'&gt; Int64Index: 1436 entries, 0 to 1435 Data columns (total 10 columns): #   Column      Non-Null Count  Dtype ---  --- 0   Price       1436 non-null   int64 1   Age         1336 non-null   float64 2   KM          1421 non-null   float64 3   FuelType    1336 non-null   object 4   HP          1430 non-null   float64 5   MetColor    1286 non-null   object 6   Automatic   1436 non-null   object 7   CC          1436 non-null   int64 8   Doors       1436 non-null   object 9   Weight      1436 non-null   int64 dtypes: float64(3), int64(3), object(4) memory usage: 123.4+ KB</pre>	<div>astype( ) method is used to explicitly convert data types from one to another</div> <div>nbytes() is used to get the total bytes consumed by the elements of the columns</div>
20	<div>Dataframe.replace( )</div>	<pre>In [10]: print('Memory when data is of category Type') car['FuelType'].astype('category').nbytes  Memory when data is of category Type  Out[10]: 1460  In [12]: #5) Cleaning column 'Doors #use of replace()  print(np.unique(car['Doors']))  car['Doors'].replace('three',3,inplace=True) car['Doors'].replace('four',4,inplace=True) car['Doors'].replace('five',5,inplace=True)  ['2' '3' '4' '5' 'five' 'four' 'three']  In [14]: car['Doors']=car['Doors'].astype('int64')  print(np.unique(car['Doors']))  [2 3 4 5]</pre>	<div>replace() is used to replace a value with the desired value</div>

21	Dataframe.isnull().sum()	<div>car.isnull().sum()</div> <div>Out[15]:</div> <table><tr><td>Price</td><td>0</td></tr><tr><td>Age</td><td>100</td></tr><tr><td>KM</td><td>15</td></tr><tr><td>FuelType</td><td>100</td></tr><tr><td>HP</td><td>6</td></tr><tr><td>MetColor</td><td>150</td></tr><tr><td>Automatic</td><td>0</td></tr><tr><td>CC</td><td>0</td></tr><tr><td>Doors</td><td>0</td></tr><tr><td>Weight</td><td>0</td></tr><tr><td colspan="2">dtype: int64</td></tr></table>	Price	0	Age	100	KM	15	FuelType	100	HP	6	MetColor	150	Automatic	0	CC	0	Doors	0	Weight	0	dtype: int64		To check the count of missing values present in each column Dataframe.isnull.sum() is used
Price	0																								
Age	100																								
KM	15																								
FuelType	100																								
HP	6																								
MetColor	150																								
Automatic	0																								
CC	0																								
Doors	0																								
Weight	0																								
dtype: int64																									