

Module 5: Data Manipulation

Case Study

edureka!

edureka!

© Brain4ce Education Solutions Pvt. Ltd.

Case Study

From the data provided on Hollywood movies:

1. Find the highest rated movie in the “Quest” story type.
2. Find the genre in which there has been the greatest number of movie releases
3. Print the names of the top five movies with the costliest budgets.
4. Is there any correspondence between the critics’ evaluation of a movie and its acceptance by the public? Find out, by plotting the net profitability of a movie against the ratings it receives on Rotten Tomatoes.

5. Perform Operations on Files

5.1: From the raw data below create a data frame

```
'first_name': ['Jason', 'Molly', 'Tina', 'Jake', 'Amy'],  
'last_name': ['Miller', 'Jacobson', ".", 'Milner', 'Cooze'],  
'age': [42, 52, 36, 24, 73],  
'preTestScore': [4, 24, 31, ".", "."],  
'postTestScore': ["25,000", "94,000", 57, 62, 70]
```

5.2: Save the dataframe into a csv file as example.csv

5.3: Read the example.csv and print the data frame

5.4: Read the example.csv without column heading

Question 5: Read the example.csv and make the index columns as 'First Name' and 'Last Name'

5.6: Print the data frame in a Boolean form as True or False. True for Null/ NaN values and false for non-null values

5.7: Read the dataframe by skipping first 3 rows and print the data frame

5.8: Load a csv file while interpreting "," in strings around numbers as thousands separators. Check the raw data 'postTestScore' column has, as thousands separator.

Comma should be ignored while reading the data. It is default behaviour, but you need to give argument to read_csv function which makes sure commas are ignored.

6. Perform Operations on Files

6.1: From the raw data below create a Pandas Series

'Amit', 'Bob', 'Kate', 'A', 'b', np.nan, 'Car', 'dog', 'cat'

- a) Print all elements in lower case
- b) Print all the elements in upper case
- c) Print the length of all the elements

6.2: From the raw data below create a Pandas Series

' Atul', 'John ', ' jack ', 'Sam'

- a) Print all elements after stripping spaces from the left and right
- b) Print all the elements after removing spaces from the left only
- c) Print all the elements after removing spaces from the right only

6.3: - Create a series from the raw data below

'India_is_big', 'Population_is_huge', np.nan, 'Has_diverse_culture'

- a) split the individual strings wherever '_' comes and create a list out of it.
- b) Access the individual element of a list
- c) Expand the elements so that all individual elements get splitted by '_' and instead of list returns individual elements

6.4: Create a series and replace either X or dog with XX-XX

'A', 'B', 'C', 'AabX', 'BacX', np.nan, 'CABA', 'dog', 'cat'

6.5: Create a series and remove dollar from the numeric values

'12', '-\$10', '\$10,000'

6.6:- Create a series and reverse all lower case words

'india 1998', 'big country', np.nan

6.7: Create pandas series and print true if value is alphanumeric in series or false if value is not alpha numeric in series.

'1', '2', '1a', '2b', '2003c'

6.8: Create pandas series and print true if value is containing 'A'

'1', '2', '1a', '2b', 'America', 'VietnAm', 'vietnam', '2003c'

6.9: Create pandas series and print in three columns value 0 or 1 is a or b or c exists in values

'a', 'a|b', np.nan, 'a|c'

6.10: Create pandas dataframe having keys and ltable and rtable as below -

'key': ['One', 'Two'], 'ltable': [1, 2]

'key': ['One', 'Two'], 'rtable': [4, 5]

Merge both the tables based of key

equiteka: