

Name	Bodhisatya Ghosh
UID no.	2021700026

Experiment 3	
HONOUR PLEDGE	<p>Please copy the following statement onto a blank piece of paper with your signature:</p> <p>“ I hereby declare that the documentation, code and output attached with this lab experiment has been completed by me in accordance with the highest standards of honesty. I confirm that I have not plagiarized OR used unauthorized materials OR given or received illegitimate help for completing this experiment. I will uphold equity and honesty in the evaluation of my work and if found guilty of plagiarism or dishonesty, will bear the consequences as outlined in the ‘integrity’ section of the lab rubrics. I am doing so in order to maintain a community built around this code of honour”</p>
PROBLEM STATEMENT :	<p>Data Integration and Reshaping:</p> <p>1. Merge two or more data frames based on a common key.</p> <ul style="list-style-type: none"> • Create a new pandas dataframe with containing 20 records and 5 attributes. • One attribute should compulsorily be a categorical variable, which is common with a categorical attribute from the CSV dataset that has been used earlier by you. • The other 4 attributes can be generated on reasonable assumptions. • Merge these 2 datasets on the common key <p>2. Concatenate multiple Data Frames vertically or horizontally</p> <ul style="list-style-type: none"> • Create 5 new rows of the same schema as the original csv dataframe. • Use a new categorical value for the common key attribute. • Concatenate this horizontally with the existing dataframe • Similarly, execute a vertical concatenation with a mock dataframe. <p>3. Pivot a Data Frame from long to wide format or vice versa</p>

- Add reasoning for using pivot. Explain with a relevant example how pivot operation is useful in data analysis

4. Stack and unstack columns or levels in the Data Frame

- Reason about the use and application of stacking and unstacking with the help of the current dataframe or another example.

5. Data Wrangling

- Experiment with other techniques in data wrangling to convert and reshape your dataframe into its final state which can be used for analysis

THEORY:

PROGRAM:

1. Merge two or more data frames based on a common key.

Create a new pandas dataframe with containing 20 records and 5 attributes:

One attribute should compulsorily be a categorical variable, which is common with a categorical attribute from the CSV dataset that has been used earlier by you.

The other 4 attributes can be generated on reasonable assumptions.

Merge these 2 datasets on the common key

```
data1 = {
    'id': range(1, 21),
    'name': ['John', 'Alice', 'Bob', 'Tom', 'Michael', 'Olivia', 'Daniel', 'Sophia', 'Charlie', 'Emma',
            'Liam', 'Mia', 'James', 'Ava', 'Matthew', 'Grace', 'William', 'Oliver', 'Hannah', 'Ivy'],
    'age': [25, 30, 22, 28, 35, 27, 31, 29, 24, 32, 26, 34, 33, 26, 36, 29, 28, 31],
    'city': ['New York', 'Los Angeles', 'Chicago', 'San Francisco', 'Miami', 'Seattle', 'Boston', 'Dallas', 'Houston', 'Minneapolis',
            'Denver', 'Portland', 'Austin', 'San Diego', 'Philadelphia', 'New York', 'Houston', 'Phoenix', 'San Antonio', 'Charlotte'],
    'salary': [7500, 9000, 6000, 8500, 11000, 9500, 8000, 6700, 7200, 9200, 8800, 6100, 9800, 7000, 10500, 7300, 8100, 6900, 9100]}

df1 = pd.DataFrame(data1)
df1.info()
```

```
[100] /> db

<class 'pandas.core.frame.DataFrame'>
Int64Index: 20 entries, 0 to 19
Data columns (total 5 columns):
 #   Column  Non-Null Count  Dtype
---  -
 0   id      20 non-null         int64
 1   name    20 non-null         object
 2   age     20 non-null         int64
 3   city    20 non-null         object
 4   salary  20 non-null         int64
dtypes: int64(4), object(1)
memory usage: 328.0 bytes
```

```
key = range(1, 21)

df2 = pd.read_csv("../rep 1/penguin_size.csv").iloc[0:20]
df2.info()
df2.info()
```

```
[107] /> db

<class 'pandas.core.frame.DataFrame'>
Int64Index: 20 entries, 0 to 19
Data columns (total 8 columns):
 #   Column  Non-Null Count  Dtype
---  -
 0   species  20 non-null      object
 1   island   20 non-null      object
 2   culmen_length_mm  19 non-null      float64
 3   culmen_depth_mm  19 non-null      float64
 4   flipper_length_mm  19 non-null      float64
 5   body_mass_g  19 non-null      float64
 6   sex      15 non-null      object
 7   id       20 non-null      int64
dtypes: float64(4), int64(1), object(3)
memory usage: 1.46 KB
```

```
merged = pd.merge(left=df1, right=df2, how='inner', on='id')

index = merged['id']

merged.drop(columns=['id'], inplace=True)
merged.set_index(index, inplace=True)

merged.info()
```

```
[108] /> db

<class 'pandas.core.frame.DataFrame'>
Int64Index: 20 entries, 1 to 20
Data columns (total 11 columns):
 #   Column  Non-Null Count  Dtype
---  -
 0   name    20 non-null      object
 1   age     20 non-null      int64
 2   city    20 non-null      object
 3   salary  20 non-null      object
 4   species  20 non-null      object
 5   island   20 non-null      object
 6   culmen_length_mm  19 non-null      float64
 7   culmen_depth_mm  19 non-null      float64
 8   flipper_length_mm  19 non-null      float64
 9   body_mass_g  19 non-null      float64
10   sex      15 non-null      object
dtypes: float64(4), int64(2), object(5)
memory usage: 1.9 KB
```

2. Concatenate multiple Data Frames vertically or horizontally

Create 5 new rows of the same schema as the original csv dataframe.

Use a new categorical value for the common key attribute.

Concatenate this horizontally with the existing dataframe.

Similarly, execute a vertical concatenation with a mock dataframe.

```
df = pd.read_csv("../penguins_size.csv")
df.tail()

100 ✓ 6/6 Python
species island culmen_length_mm culmen_depth_mm flipper_length_mm body_mass_g sex
339 Gentoo Biscoe NaN NaN NaN NaN
340 Gentoo Biscoe 46.8 14.3 215.0 4850.0 FEMALE
341 Gentoo Biscoe 50.4 15.7 222.0 5750.0 MALE
342 Gentoo Biscoe 45.2 14.8 212.0 5200.0 FEMALE
343 Gentoo Biscoe 49.3 16.1 213.0 5400.0 MALE

# Create a new dataframe for horizontal concatenation
new_data_horizontal = {
    'species': ['New_Species_1', 'New_Species_2', 'New_Species_3', 'New_Species_4', 'New_Species_5'],
    'island': ['New_Island'] * 5,
    'culmen_length_mm': [41.2, 38.5, 40.4, 35.5, 37.8],
    'culmen_depth_mm': [14.5, 12.2, 17.0, 16.5, 19.8],
    'flipper_length_mm': [185, 180, 189, 195, 200],
    'body_mass_g': [1200, 1500, 3500, 7000, 4000],
    'sex': ['MALE', 'FEMALE', 'FEMALE', 'MALE', 'FEMALE']
}

df_horizontal = pd.DataFrame(new_data_horizontal)

# Horizontal Concatenation
df_concat_horizontal = pd.concat([df, df_horizontal], axis=1)
df_concat_horizontal.tail()

100 ✓ 6/6 Python
species island culmen_length_mm culmen_depth_mm flipper_length_mm body_mass_g sex species island culmen_length_mm culmen_depth_mm flipper_length_mm body_mass_g sex
339 Gentoo Biscoe NaN NaN NaN NaN NaN NaN
340 Gentoo Biscoe 46.8 14.3 215.0 4850.0 FEMALE NaN NaN NaN NaN NaN NaN
341 Gentoo Biscoe 50.4 15.7 222.0 5750.0 MALE NaN NaN NaN NaN NaN NaN
342 Gentoo Biscoe 45.2 14.8 212.0 5200.0 FEMALE NaN NaN NaN NaN NaN NaN
343 Gentoo Biscoe 49.3 16.1 213.0 5400.0 MALE NaN NaN NaN NaN NaN NaN
```

```
# Create a new dataframe for vertical concatenation
new_data_vertical = {
    'species': ['New_Species_5', 'New_Species_7', 'New_Species_8', 'New_Species_9', 'New_Species_10'],
    'island': ['New_Island'] * 5,
    'culmen_length_mm': [42.0, 39.0, 38.5, 41.0, 40.5],
    'culmen_depth_mm': [17.0, 16.8, 17.5, 16.2, 18.8],
    'flipper_length_mm': [190, 187, 180, 200, 195],
    'body_mass_g': [1800, 1600, 3500, 4000, 3000],
    'sex': ['FEMALE', 'MALE', 'FEMALE', 'MALE', 'FEMALE']
}

df_vertical = pd.DataFrame(new_data_vertical)

# Vertical Concatenation
df_concat_vertical = pd.concat([df, df_vertical], ignore_index=True)
df_concat_vertical.tail()

100 ✓ 6/6 Python
species island culmen_length_mm culmen_depth_mm flipper_length_mm body_mass_g sex
344 New_Species_5 New_Island 42.0 17.0 192.0 3800.0 FEMALE
345 New_Species_7 New_Island 39.0 16.0 192.0 3800.0 MALE
346 New_Species_8 New_Island 38.5 17.5 180.0 3500.0 FEMALE
347 New_Species_9 New_Island 41.0 16.2 200.0 4000.0 MALE
348 New_Species_10 New_Island 40.5 18.0 195.0 3000.0 FEMALE
```

3. Pivot a Data Frame from long to wide format or vice versa

Add reasoning for using pivot. Explain with a relevant example how pivot operation is useful in data analysis

```
# Create a random meaningful dataframe
np.random.seed(42)

data = {
    'Date': pd.date_range(start='2022-01-01', periods=12, freq='W'),
    'Category': np.random.choice(['A', 'B', 'C'], size=12),
    'Value': np.random.randint(1, 100, size=12),
    'Location': np.random.choice(['X', 'Y'], size=12)
}

df = pd.DataFrame(data)

# Display the original dataframe
print("Original DataFrame")
print(df)

100 ✓ 6/6 Python
Original DataFrame:
   Date Category Value Location
0 2022-01-01      C    88      X
1 2022-01-08      A    74      X
2 2022-01-15      C     3      Y
3 2022-01-22      C    23      Y
4 2022-01-29      A    51      Y
5 2022-02-05      A     2      X
6 2022-02-12      C    69      Y
7 2022-02-19      B    38      X
8 2022-02-26      C     3      X
9 2022-03-05      C     2      X
10 2022-11-30      C    64      X
11 2022-12-31      C    69      X

# Pivot from long to wide format
df_pivot_wide = df.pivot_table(df, values='Value', index='Date', columns=['Location', 'Category'])

# Display the dataframe after pivoting to wide format
print("DataFrame after pivoting to wide format")
print(df_pivot_wide)

100 ✓ 6/6 Python
DataFrame after Pivoting to Wide Format:
   Date      Value
   Location
Category A B C A C
Date
2022-01-01 NaN NaN 88.0 NaN NaN
2022-01-08 NaN NaN NaN NaN NaN
2022-01-15 NaN NaN NaN NaN 3.0
2022-01-22 NaN NaN NaN NaN 23.0
2022-01-29 NaN NaN 51.0 NaN
2022-02-05 NaN NaN NaN NaN 2.0
2022-02-12 NaN 69.0 NaN NaN NaN
2022-02-19 NaN 38.0 NaN NaN NaN
2022-02-26 NaN NaN 3.0 NaN NaN
2022-11-30 NaN NaN 64.0 NaN NaN
2022-12-31 NaN NaN 69.0 NaN NaN
```

Data pivoting enables you to rearrange the columns and rows in a report so you can view data from different perspectives. If you pivot the objects on the report, so that the objects that were in the columns are now in the rows, and the objects that were in the rows are now in the columns, much of the data is easier to read and compare.

4. Stack and unstack columns or levels in the Data Frame

Reason about the use and application of stacking and unstacking with the help of the current dataframe or another example.

```
# Pivot from wide to long format
df_pivot_long = df_pivot_wide.stack(level=['Location', 'Category'])

# Display the dataframe after pivoting to long format
print("DataFrame after Pivoting to Long format")
print(df_pivot_long)

100 ✓ 6/6 Python
DataFrame after Pivoting to Long Format:
   Date Location Category Value
0 2022-01-01 X C 88.0
1 2022-01-08 X A 74.0
2 2022-01-15 Y C 3.0
3 2022-01-22 Y A 23.0
4 2022-01-29 X A 51.0
5 2022-02-05 X A 2.0
6 2022-02-12 Y C 69.0
7 2022-02-19 X B 38.0
8 2022-02-26 X C 3.0
9 2022-03-05 X C 2.0
```

5. Data Wrangling

Experiment with other techniques in data wrangling to convert and reshape your dataframe into its final state which can be used for analysis

```
# Original Dataframe
print("Original Dataframe:")
print(df)

# Data Wrangling Techniques

# 1. Handling Missing Values
df_cleaned = df.dropna()
print("Dataframe after Handling Missing Values:")
print(df_cleaned)
```

[14] ✓ 6/6

Original Dataframe:

	Date	Category	Value	Location
0	2022-01-31	C	88	X
1	2022-02-28	A	24	X
2	2022-03-31	C	3	Y
3	2022-04-30	C	22	Y
4	2022-05-31	A	53	Y
5	2022-06-30	A	2	X
6	2022-07-31	C	88	Y
7	2022-08-31	B	30	X
8	2022-09-30	C	38	X
9	2022-10-31	C	2	X
10	2022-11-30	C	64	X
11	2022-12-31	C	60	X

Dataframe after Handling Missing Values:

	Date	Category	Value	Location
0	2022-01-31	C	88	X
1	2022-02-28	A	24	X
2	2022-03-31	C	3	Y
3	2022-04-30	C	22	Y
4	2022-05-31	A	53	Y
5	2022-06-30	A	2	X
6	2022-07-31	C	88	Y
7	2022-08-31	B	30	X
8	2022-09-30	C	38	X
9	2022-10-31	C	2	X
10	2022-11-30	C	64	X
11	2022-12-31	C	60	X

```
# 2. Grouping and Aggregation
df_grouped = df.groupby(['Category', 'Location']).agg({'Value': 'mean'}).reset_index()
print("Dataframe after Grouping and Aggregation:")
print(df_grouped)
```

[15] ✓ 6/6

Dataframe after Grouping and Aggregation:

	Category	Location	Value
0		X	13.000000
1	A	Y	53.000000
2	B	X	30.000000
3	C	X	56.400000
4	C	Y	37.466667

```
# 3. Merging Dataframes
additional_data = {
    "Category": ['A', 'B', 'C'],
    "Additional_Info": ['Info_A', 'Info_B', 'Info_C']
}
df_additional_info = pd.DataFrame(additional_data)

df_merged = pd.merge(df_grouped, df_additional_info, on='Category', how='left')
print("Dataframe after Merging with Additional Information:")
print(df_merged)
```

[16] ✓ 6/6

Dataframe after Merging with Additional Information:

	Category	Location	Value	Additional_Info
0	A	X	13.000000	Info_A
1	A	Y	53.000000	Info_A
2	B	X	30.000000	Info_B
3	C	X	56.400000	Info_C
4	C	Y	37.466667	Info_C

RESULT:

References:	pandas.merge — pandas 2.2.0 documentation (pydata.org) Pivoting data (microstrategy.com) Reshaping and pivot tables — pandas 2.2.0 documentation (pydata.org)
CONCLUSION: In conclusion, the above used data manipulation techniques, including merging, concatenating, pivoting, stacking, and data wrangling have been performed successfully	