# LESSON 4.2

## SECTION-1: loc and iloc

```python
import pandas as pd
```

```python
file = "Resources/sampleData.csv"
```

```python
original_df = pd.read_csv(file)
original_df.head()
```

```python
# Set new index to last_name
df = original_df.set_index("last_name")
df.head()
```

```python
# Grab the data contained within the "Berry" row and the "Phone Number" column
berry_phone = df.loc["Berry", "Phone Number"]
print("Using Loc: " + berry_phone)

also_berry_phone = df.iloc[1, 2]
print("Using Iloc: " + also_berry_phone)
```

```python
# Grab the first five rows of data and the columns from "id" to "Phone Number"
# The problem with using "last_name" as the index is that the values are not unique so duplicates are returned
# If there are duplicates and loc[] is being used, Pandas will return an error
richardson_to_morales = df.loc[["Richardson", "Berry", "Hudson",
                                "Mcdonald", "Morales"], ["id", "first_name", "Phone Number"]]
print(richardson_to_morales)

print()

# Using iloc[] will not find duplicates since a numeric index is always unique
also_richardson_to_morales = df.iloc[0:4, 0:3]
print(also_richardson_to_morales)
```

```python
# The following will select all rows for columns `first_name` and `Phone Number`
df.loc[:, ["first_name", "Phone Number"]].head()
```

```python
# the following logic test/conditional statement returns a series of boolean values
named_billy = df["first_name"] == "Billy"
named_billy.head()
```

```python
# Loc and Iloc also allow for conditional statments to filter rows of data
# using Loc on the Logic test above only returns rows where the result is True
only_billys = df.loc[df["first_name"] == "Billy", :]
print(only_billys)

print()

# Multiple conditions can be set to narrow down or widen the filter
only_billy_and_peter = df.loc[(df["first_name"] == "Billy") | (
    df["first_name"] == "Peter"), :]
print(only_billy_and_peter)
```

# SECTION-2: DATA CLEANING

```
In [ ]:    1  # Dependencies
           2  import pandas as pd
```

```
In [ ]:    1  # Name of the CSV file
           2  file = 'Resources/donors2008.csv'
```

```
In [ ]:    1  # The correct encoding must be used to read the CSV in pandas
           2  df = pd.read_csv(file, encoding="ISO-8859-1")
```

```
In [ ]:    1  # Preview of the DataFrame
           2  # Note that FIELD8 is likely a meaningless column
           3  df.head()
```

```
In [ ]:    1  # Delete extraneous column
           2  del df['FIELD8']
           3  df.head()
```

## Pandas DataFrame: drop() function

As an alternative we can use the drop() function as well.

```
In [ ]:    1  # Delete a column using drop()
           2  #df_alt = df.drop(["FIELD8"], axis=1)
           3  #df_alt.head()
```

```
In [ ]:    1  # Identify incomplete rows
           2  df.count()
```

```
In [ ]:    1  # Drop all rows with missing information
           2  df = df.dropna(how='any')
```

```
In [ ]:    1  # Verify dropped rows
           2  df.count()
```

```
In [ ]:    1  # The Amount column is the wrong data type. It should be numeric.
           2  df.dtypes
```

```
In [ ]:    1  # Use pd.to_numeric() method to convert the datatype of the Amount column
           2  df['Amount'] = pd.to_numeric(df['Amount'])
```

```
In [ ]:    1  # Verify that the Amount column datatype has been made numeric
           2  df['Amount'].dtype
```

```
In [ ]:    1  # Display an overview of the Employers column
           2  df['Employer'].value_counts()
```

```
In [ ]:    1  # Clean up Employer category. Replace 'Self Employed' and 'Self' with 'Self-Employed'
           2  df['Employer'] = df['Employer'].replace(
           3      {'Self Employed': 'Self-Employed', 'Self': 'Self-Employed'})
```

```
In [ ]:    1  # Verify clean-up.
           2  df['Employer'].value_counts()
```

```
In [ ]:    1  df['Employer'] = df['Employer'].replace({'Not Employed': 'Unemployed'})
           2  df['Employer'].value_counts()
```

```
In [ ]:    1  # Display a statistical overview
           2  # We can infer the maximum allowable individual contribution from 'max'
           3  df.describe()
```

**PANDAS RECAP**

```python
1  # Import the Pandas library
2  import pandas as pd
```

```python
1  # Create a reference the CSV file desired
2  csv_path = "Resources/ufoSightings.csv"
3
4  # Read the CSV into a Pandas DataFrame
5  ufo_df = pd.read_csv(csv_path)
6
7  # Print the first five rows of data to the screen
8  ufo_df.head()
```

```python
1  # Check to see if there are any rows with missing data
2  ufo_df.count()
```

```python
1  # Remove the rows with missing data
2  clean_ufo_df = ufo_df.dropna(how="any")
3  clean_ufo_df.count()
```

```python
1   # Collect a list of sightings seen in the US
2   columns = [
3       "datetime",
4       "city",
5       "state",
6       "country",
7       "shape",
8       "duration (seconds)",
9       "duration (hours/min)",
10      "comments",
11      "date posted"
12  ]
13
14  # Filter the data so that only those sightings in the US are in a DataFrame
15  usa_ufo_df = clean_ufo_df.loc[clean_ufo_df["country"] == "us", columns]
16  usa_ufo_df.head()
```

```python
1  # Count how many sightings have occured within each state
2  state_counts = usa_ufo_df["state"].value_counts()
3  state_counts
```

```python
1  # Convert the state_counts Series into a DataFrame
2  state_ufo_counts_df = pd.DataFrame(state_counts)
3  state_ufo_counts_df.head()
```

```python
1  # Convert the column name into "Sum of Sightings"
2  state_ufo_counts_df = state_ufo_counts_df.rename(
3      columns={"state": "Sum of Sightings"})
4  state_ufo_counts_df.head()
```

```python
1  # Want to add up the seconds UFOs are seen? There is a problem
2  # Problem can be seen by examining datatypes within the DataFrame
3  usa_ufo_df.dtypes
```

```python
1  # Using astype() to convert a column's data into floats
2  usa_ufo_df.loc[:, "duration (seconds)"] = usa_ufo_df["duration (seconds)"].astype("float")
3  usa_ufo_df.dtypes
```

```python
1  # Now it is possible to find the sum of seconds
2  usa_ufo_df["duration (seconds)"].sum()
```

## SECTION-3: GROUPBY

```python
1  # Import Dependencies
2  import pandas as pd
```

```python
1  # Create a reference the CSV file desired
2  csv_path = "Resources/ufoSightings.csv"
3
4  # Read the CSV into a Pandas DataFrame
5  ufo_df = pd.read_csv(csv_path)
6
7  # Print the first five rows of data to the screen
8  ufo_df.head()
```

```python
1  # Remove the rows with missing data
2  clean_ufo_df = ufo_df.dropna(how="any")
3  clean_ufo_df.count()
```

```python
1  clean_ufo_df.head()
```

```python
1  # Converting the "duration (seconds)" column's values to numeric
2  converted_ufo_df = clean_ufo_df.copy()
3  converted_ufo_df["duration (seconds)"] = converted_ufo_df.loc[:, "duration (seconds)"].astype(float)
```

```python
1  converted_ufo_df.head()
```

```python
1  # Filter the data so that only those sightings in the US are in a DataFrame
2  usa_ufo_df = converted_ufo_df.loc[converted_ufo_df["country"] == "us", :]
3  usa_ufo_df.head()
```

```python
1  # Count how many sightings have occured within each state
2  state_counts = usa_ufo_df["state"].value_counts()
3  state_counts.head()
```

```python
1  # Using GroupBy in order to separate the data into fields according to "state" values
2  grouped_usa_df = usa_ufo_df.groupby(['state'])
3
4  # The object returned is a "GroupBy" object and cannot be viewed normally...
5  print(grouped_usa_df)
6
7  # In order to be visualized, a data function must be used...
8  grouped_usa_df.count().head(10)
```

```python
1  grouped_usa_df["duration (seconds)"].sum()
```

```python
1  # Since "duration (seconds)" was converted to a numeric time, it can now be summed up per state
2  state_duration = grouped_usa_df["duration (seconds)"].sum()
3  state_duration.head()
```

```python
1  # Creating a new DataFrame using both duration and count
2  state_summary_df = pd.DataFrame({"Number of Sightings": state_counts,
3                                   "Total Visit Time": state_duration})
4  state_summary_df.head()
```

```python
1  # It is also possible to group a DataFrame by multiple columns
2  # This returns an object with multiple indexes, however, which can be harder to deal with
3  grouped_international_data = converted_ufo_df.groupby(['country', 'state'])
4
5  grouped_international_data.count().head(20)
```

```python
1  # Converting a GroupBy object into a DataFrame
2  international_duration_df = pd.DataFrame(
3      grouped_international_data["duration (seconds)"].sum())
4  international_duration_df.head(10)
```

**SECTION 4: SORTING**

```python
1  # Import Dependencies
2  import pandas as pd
```

```python
1  csv_path = "Resources/Happiness_2017.csv"
2  happiness_df = pd.read_csv(csv_path)
3  happiness_df.head()
```

```python
1  # Sorting the DataFrame based on "Freedom" column
2  # Will sort from lowest to highest if no other parameter is passed
3  freedom_df = happiness_df.sort_values("Freedom")
4  freedom_df.head()
```

```python
1  # To sort from highest to lowest, ascending=False must be passed in
2  freedom_df = happiness_df.sort_values("Freedom", ascending=False)
3  freedom_df.head()
```

```python
1  # It is possible to sort based upon multiple columns
2  family_and_generosity_df = happiness_df.sort_values(
3      ["Family", "Generosity"], ascending=False)
4  family_and_generosity_df.head()
```

```python
1  # The index can be reset to provide index numbers based on the new rankings.
2  new_index_df = family_and_generosity_df.reset_index(drop=True)
3  new_index_df.head()
```