

DATA WRANGLING

In [187]:

```
1 import pandas as pd
2
3 #online link
4 #url='https://raw.githubusercontent.com/chrisalbon/simulated_datasets/master/titanic.csv'
5
6 #csv file
7 url = 'C:/Users/Prerna/Desktop/ML_jupyter_notenooks/titanic.csv'
8
9 #Creating a dataframe
10 dataframe = pd.read_csv(url)
11
12 #printing the dataframe
13 dataframe
14
```

Out[187]:

	Name	PClass	Age	Sex	Survived	SexCode
0	Allen, Miss Elisabeth Walton	1st	29.00	female	1	1
1	Allison, Miss Helen Loraine	1st	2.00	female	0	1
2	Allison, Mr Hudson Joshua Creighton	1st	30.00	male	0	0
3	Allison, Mrs Hudson JC (Bessie Waldo Daniels)	1st	25.00	female	0	1
4	Allison, Master Hudson Trevor	1st	0.92	male	1	0
5	Anderson, Mr Harry	1st	47.00	male	1	0
6	Andrews, Miss Kornelia Theodosia	1st	63.00	female	1	1
7	Andrews, Mr Thomas, jr	1st	39.00	male	0	0
8	Appleton, Mrs Edward Dale (Charlotte Lamson)	1st	58.00	female	1	1
9	Artagaveytia, Mr Ramon	1st	71.00	male	0	0
10	Astor, Colonel John Jacob	1st	47.00	male	0	0

```
In [188]: 1 #Reading first five records with head()
          2 dataframe.head(5)
```

Out[188]:

	Name	PClass	Age	Sex	Survived	SexCode
0	Allen, Miss Elisabeth Walton	1st	29.00	female	1	1
1	Allison, Miss Helen Loraine	1st	2.00	female	0	1
2	Allison, Mr Hudson Joshua Creighton	1st	30.00	male	0	0
3	Allison, Mrs Hudson JC (Bessie Waldo Daniels)	1st	25.00	female	0	1
4	Allison, Master Hudson Trevor	1st	0.92	male	1	0

```
In [189]: 1 #Reading last five records with tail()
          2 dataframe.tail(5)
```

Out[189]:

	Name	PClass	Age	Sex	Survived	SexCode
1308	Zakarian, Mr Artun	3rd	27.0	male	0	0
1309	Zakarian, Mr Maprieder	3rd	26.0	male	0	0
1310	Zenni, Mr Philip	3rd	22.0	male	0	0
1311	Lievens, Mr Rene	3rd	24.0	male	0	0
1312	Zimmerman, Leo	3rd	29.0	male	0	0

```
In [190]: 1 #showing dimensions of the data
          2 dataframe.shape
```

Out[190]: (1313, 6)

(1313,6) = (no of rows, no of columns)

```
In [191]: 1 #get the descriptive statistics of data
          2 dataframe.describe()
```

Out[191]:

	Age	Survived	SexCode
count	756.000000	1313.000000	1313.000000
mean	30.397989	0.342727	0.351866
std	14.259049	0.474802	0.477734
min	0.170000	0.000000	0.000000
25%	21.000000	0.000000	0.000000
50%	28.000000	0.000000	0.000000
75%	39.000000	1.000000	1.000000
max	71.000000	1.000000	1.000000

Navigating the dataframes using loc and iloc.

All rows in a pandas DataFrame have a unique index value.

By default, this index is an integer indicating the row position in the DataFrame however, it does not have to be.

DataFrame indexes can be set to be unique alphanumeric strings or customer numbers.

To select individual rows and slices of rows, pandas provides two methods:

- loc is useful when the index of the DataFrame is a label (e.g., a string).
- iloc works by looking for the position in the DataFrame.

```
In [192]: 1 # selecting one of more rows with iloc.  
2 #Selecting the first row  
3 dataframe.iloc[0]
```

```
Out[192]: Name      Allen, Miss Elisabeth Walton  
PClass      1st  
Age         29  
Sex         female  
Survived    1  
SexCode     1  
Name: 0, dtype: object
```

```
In [193]: 1 #Selecting more than one rows using slicing with iloc  
2 #selecting first 4 rows from the dataset  
3 dataframe.iloc[1:4]
```

```
Out[193]:
```

	Name	PClass	Age	Sex	Survived	SexCode
1	Allison, Miss Helen Loraine	1st	2.0	female	0	1
2	Allison, Mr Hudson Joshua Creighton	1st	30.0	male	0	0
3	Allison, Mrs Hudson JC (Bessie Waldo Daniels)	1st	25.0	female	0	1

```
In [194]: 1 #selecting first 4 rows from the dataset  
2 dataframe.iloc[:4]  
3  
4 #one more variation for selecting first 4 rows  
5 #dataframe.iloc[0:4]
```

```
Out[194]:
```

	Name	PClass	Age	Sex	Survived	SexCode
0	Allen, Miss Elisabeth Walton	1st	29.0	female	1	1
1	Allison, Miss Helen Loraine	1st	2.0	female	0	1
2	Allison, Mr Hudson Joshua Creighton	1st	30.0	male	0	0
3	Allison, Mrs Hudson JC (Bessie Waldo Daniels)	1st	25.0	female	0	1

Using loc for selecting the rows.

DataFrames do not need to be numerically indexed.

We can set the index of a Data- Frame to any value where the value is unique to each row.

For example, we can set the index to be passenger names and then select rows using a name.

Use `set_index()` for setting the index.

```
In [195]: 1 #Set the Name column as index value for dataframe
          2 dataframe = dataframe.set_index(dataframe['Name'])
          3
          4 #Select the rows with Loc.
          5 dataframe.loc['Allison, Miss Helen Loraine']
```

```
Out[195]: Name      Allison, Miss Helen Loraine
          PCClass      1st
          Age          2
          Sex          female
          Survived      0
          SexCode      1
          Name: Allison, Miss Helen Loraine, dtype: object
```

```
In [196]: 1 #Selecting the rows with Loc using slicing
          2 dataframe.loc['Allen, Miss Elisabeth Walton':'Allison, Mr Hudson Joshua Creighton']
```

Out[196]:

	Name	PClass	Age	Sex	Survived	SexCode	
	Allen, Miss Elisabeth Walton	Allen, Miss Elisabeth Walton	1st	29.0	female	1	1
	Allison, Miss Helen Loraine	Allison, Miss Helen Loraine	1st	2.0	female	0	1
	Allison, Mr Hudson Joshua Creighton	Allison, Mr Hudson Joshua Creighton	1st	30.0	male	0	0

Selecting Rows Based on Conditionals

In [197]:

```
1 import pandas as pd
2
3 #Creating URL
4 url = 'C:/Users/Prerna/Desktop/ML_jupyter_notenooks/titanic.csv'
5
6 #Creating a dataframe
7 dataframe = pd.read_csv(url)
8
9 # 1) Select all the women on the titanic and display top 5 records of women
10
11 dataframe[dataframe['Sex']=='female'].head(5)
12
```

Out[197]:

	Name	PClass	Age	Sex	Survived	SexCode
0	Allen, Miss Elisabeth Walton	1st	29.0	female	1	1
1	Allison, Miss Helen Loraine	1st	2.0	female	0	1
3	Allison, Mrs Hudson JC (Bessie Waldo Daniels)	1st	25.0	female	0	1
6	Andrews, Miss Kornelia Theodosia	1st	63.0	female	1	1
8	Appleton, Mrs Edward Dale (Charlotte Lamson)	1st	58.0	female	1	1

In [198]:

```
1 # Filter rows with multiple conditions
2 '''2) select all the 2nd class male passengers on the titanic having
3 age more than 50 years '''
4
5 dataframe[(dataframe['Sex']=='male') & (dataframe['Age']>50)
6           & (dataframe['PClass']=='2nd')]
7
```

Out[198]:

	Name	PClass	Age	Sex	Survived	SexCode
328	Ashby, Mr John	2nd	57.0	male	0	0
333	Bateman, Rev Robert James	2nd	51.0	male	0	0
362	Carter, Rev Ernest Courtenay	2nd	54.0	male	0	0
364	Chapman, Mr Charles Henry	2nd	52.0	male	0	0
447	Hodges, Mr Henry Price	2nd	52.0	male	0	0
505	Mitchell, Mr Henry Michael	2nd	71.0	male	0	0
506	Moraweck, Dr Ernest	2nd	54.0	male	0	0
509	Myles, Mr Thomas Francis	2nd	64.0	male	0	0
556	Sjostedt, Mr Ernst Adolf	2nd	59.0	male	0	0

In [200]:

```
1 #3) select 1st class male passengers that has survived or more than 50 years in titanic
2
3 dataframe[(((dataframe['Sex']=='male') & (dataframe['PClass']=='1st'))
4           & ((dataframe['Survived']==1) | (dataframe['Age']>50)))]
```

Out[200]:

	Name	PClass	Age	Sex	Survived	SexCode
4	Allison, Master Hudson Trevor	1st	0.92	male	1	0
5	Anderson, Mr Harry	1st	47.00	male	1	0
9	Artagaveytia, Mr Ramon	1st	71.00	male	0	0
13	Barkworth, Mr Algernon H	1st	NaN	male	1	0
18	Beckwith, Mr Richard Leonard	1st	37.00	male	1	0
20	Behr, Mr Karl Howell	1st	26.00	male	1	0
22	Bishop, Mr Dickinson H	1st	25.00	male	1	0
24	Bjornstrm-Steffansson, Mr Mauritz Hakan	1st	28.00	male	1	0
26	Blank, Mr Henry	1st	39.00	male	1	0
32	Bradley, Mr George	1st	NaN	male	1	0
40	Calderhead, Mr Edward P	1st	NaN	male	1	0
43	Cardeza, Mr Thomas Drake Martinez	1st	36.00	male	1	0
47	Carter, Mr William Ernest	1st	36.00	male	1	0
50	Carter, Master William T II	1st	11.00	male	1	0
57	Chambers, Mr Norman Campbell	1st	27.00	male	1	0
60	Chevre, Mr Paul	1st	NaN	male	1	0
72	Crosby, Captain Edward Gifford	1st	70.00	male	0	0
77	Daly, Mr Peter Denis	1st	NaN	male	1	0
78	Daniel, Mr Robert Williams	1st	27.00	male	1	0
82	Dick, Mr Albert Adrian	1st	31.00	male	1	0
84	Dodge, Dr Washington	1st	NaN	male	1	0
86	Dodge, Master Washington	1st	4.00	male	1	0

	Name	PClass	Age	Sex	Survived	SexCode
90	Duff Gordon, Sir Cosmo Edmund	1st	49.00	male	1	0
97	Flynn, Mr John Irving	1st	NaN	male	1	0
103	Fortune, Mr Mark	1st	64.00	male	0	0
106	Frauenthal, Dr Henry William	1st	49.00	male	1	0
108	Frauenthal, Mr Isaac Gerald	1st	44.00	male	1	0
110	Frolicher-Stehli, Mr Maxmillian	1st	60.00	male	1	0
117	Goldenberg, Mr Samuel L	1st	49.00	male	1	0
119	Goldschmidt, Mr George B	1st	71.00	male	0	0
...
221	Ryerson, Master John Borie	1st	13.00	male	1	0
223	Saalfeld, Mr Adolphe	1st	NaN	male	1	0
224	Salomon, Mr Abraham L	1st	NaN	male	1	0
226	Seward, Mr Frederic Kimber	1st	34.00	male	1	0
228	Silverthorne, Mr Spencer Victor	1st	36.00	male	1	0
231	Simonius-Blumer, Col Alfons	1st	56.00	male	1	0
232	Sloper, Mr William Thompson	1st	28.00	male	1	0
233	Smart, Mr John Montgomery	1st	56.00	male	0	0
234	Smith, Mr James Clinch	1st	56.00	male	0	0
238	Snyder, Mr John Pillsbury	1st	24.00	male	1	0
240	Spedden, Mr Frederick Oakley	1st	45.00	male	1	0
242	Spedden, Master Robert Douglas	1st	6.00	male	1	0
243	Spencer, Mr William Augustus	1st	57.00	male	0	0
245	Staehlin, Dr Max	1st	32.00	male	1	0
246	Stead, Mr William Thomas	1st	62.00	male	0	0
247	Stengel, Mr Charles Emil Henry	1st	54.00	male	1	0
252	Straus, Mr Isidor	1st	67.00	male	0	0

	Name	PClass	Age	Sex	Survived	SexCode
254	Sutton, Mr Frederick	1st	61.00	male	0	0
256	Taussig, Mr Emil	1st	52.00	male	0	0
259	Taylor, Mr Elmer Zebley	1st	48.00	male	1	0
263	Thayer, Mr John Borland, jr	1st	17.00	male	1	0
266	Tucker, Mr Gilbert Milligan, jr	1st	31.00	male	1	0
268	Van Derhoef, Mr Wyckoff	1st	61.00	male	0	0
270	Warren, Mr Frank Manley	1st	64.00	male	0	0
272	Weir, Col John	1st	60.00	male	0	0
274	White, Mr Percival Wayland	1st	54.00	male	0	0
276	Wick, Mr George Dennick	1st	57.00	male	0	0
283	Williams, Mr Charles Duane	1st	51.00	male	0	0
285	Williams, Mr Richard Norris II	1st	21.00	male	1	0
286	Woolner, Mr Hugh	1st	NaN	male	1	0

86 rows × 6 columns

Replacing Values

You can replace the values in the columns with `replace()`.

Syntax: `dataframe(columnname).replace(value to be replaced, Value to replace with)`

```
In [201]: 1 import pandas as pd
2 #Creating URL
3 url = 'C:/Users/Prerna/Desktop/ML_jupyter_notenooks/titanic.csv'
4
5 #Creating a dataframe
6 dataframe = pd.read_csv(url)
7
8 #Replacing values in Survived column 1 with Yes.
9 dataframe['Survived'].replace(1,'yes').head(5)
10
```

```
Out[201]: 0    yes
1         0
2         0
3         0
4    yes
Name: Survived, dtype: object
```

```
In [202]: 1 #Replacing values of 1 with Yes and 0 with NO, and Pclass in whole dataframe.
2 new_df = dataframe.replace([0,1,'1st','2nd','3rd'],
3                             ['Yes','No','First','Second','Third'])
```

Renaming Columns and Deleting Columns

In [203]:

```
1
2 #Drop column with drop() with axis=1 for column
3 new_df = new_df.drop('SexCode',axis=1)
4
5 #Rename column PClass to Passenger class using rename()
6 new_df = new_df.rename(columns = {'PClass': 'Passenger Class', 'Sex': 'Gender'})
7 new_df.head(5)
```

Out[203]:

	Name	Passenger Class	Age	Gender	Survived
0	Allen, Miss Elisabeth Walton	First	29	female	No
1	Allison, Miss Helen Loraine	First	2	female	Yes
2	Allison, Mr Hudson Joshua Creighton	First	30	male	Yes
3	Allison, Mrs Hudson JC (Bessie Waldo Daniels)	First	25	female	Yes
4	Allison, Master Hudson Trevor	First	0.92	male	No

In [204]:

```
1  #Merging the whole code
2
3  import pandas as pd
4  #Creating URL
5  url = 'C:/Users/Prerna/Desktop/ML_jupyter_notenooks/titanic.csv'
6
7  #Creating a dataframe
8  dataframe = pd.read_csv(url)
9
10 #Drop column with drop() with axis=1 for column
11 dataframe = dataframe.drop('SexCode',axis=1)
12
13 #Rename column PClass to Passenger class using rename()
14 dataframe = dataframe.rename(columns = {'PClass': 'Passenger Class',
15                                         'Sex': 'Gender'})
16
17 #Replacing values of 1 with Yes and 0 with NO, and Pclass in whole dataframe.
18 dataframe = dataframe.replace([0,1,'1st','2nd','3rd'],
19                               ['Yes','No','First','Second','Third'])
20
21 dataframe.head(5)
22
```

Out[204]:

	Name	Passenger Class	Age	Gender	Survived
0	Allen, Miss Elisabeth Walton	First	29	female	No
1	Allison, Miss Helen Loraine	First	2	female	Yes
2	Allison, Mr Hudson Joshua Creighton	First	30	male	Yes
3	Allison, Mrs Hudson JC (Bessie Waldo Daniels)	First	25	female	Yes
4	Allison, Master Hudson Trevor	First	0.92	male	No

Finding the Minimum, Maximum, Sum, Average, and Count

In [205]:

```
1  #Creating URL
2  import pandas as pd
3  url = 'C:/Users/Prerna/Desktop/ML_jupyter_notenooks/titanic.csv'
4
5  #Creating a dataframe
6  dataframe = pd.read_csv(url)
7
8  #Printing the maximum with max() of age
9  print("Maximum Age:", dataframe['Age'].max())
10
11 #Printing the minimum with min() of age
12 print("Minimum Age:", dataframe['Age'].min())
13
14 #Printing the Sum with sum() of age
15 print("Total Age:", dataframe['Age'].sum())
16
17 #Printing the mean with mean() of age
18 print("Mean of Age:", dataframe['Age'].mean())
19
20 #Printing the median with median() of age
21 print("Median of Age:", dataframe['Age'].median())
22
23 #Printing the mode with mode() of age
24 print("Mode of Age:", dataframe['Age'].mode())
25
26 #Printing the count with count() of age
27 print("Count of Age:", dataframe['Age'].count())
28
29 #Printing the variance with var() of age
30 print("Variance of Age:", dataframe['Age'].var())
31
32 #Printing the standard deviation with std() of age
33 print("Standard Deviation of Age:", dataframe['Age'].std())
```

Maximum Age: 71.0

Minimum Age: 0.17

Total Age: 22980.88

Mean of Age: 30.397989417989415

Median of Age: 28.0

Mode of Age: 0 22.0

```
dtype: float64
Count of Age: 756
Variance of Age: 203.32047012439133
Standard Deviation of Age: 14.259048710359023
```

```
In [206]: 1 #Show counts
          2 dataframe.count()
```

```
Out[206]: Name      1313
          PClass    1313
          Age       756
          Sex       1313
          Survived  1313
          SexCode   1313
          dtype: int64
```

Finding Unique Values

```
In [207]: 1 #To find the unique values in a column with unique()
          2
          3 #Creating URL
          4 url = 'C:/Users/Prerna/Desktop/ML_jupyter_notenooks/titanic.csv'
          5
          6 #Creating a dataframe
          7 dataframe = pd.read_csv(url)
          8
          9 #select unique values
         10 dataframe['PClass'].unique()
         11
```

```
Out[207]: array(['1st', '2nd', '*', '3rd'], dtype=object)
```

```
In [208]: 1 '''value_counts() will display all unique values with the number
          2 of times each value appears '''
          3
          4 dataframe['PClass'].value_counts()
```

```
Out[208]: 3rd    711
          1st    322
          2nd    279
          *         1
          Name: PClass, dtype: int64
```

```
In [209]: 1 # Show number of unique values with nunique()
          2
          3 dataframe['PClass'].nunique()
```

```
Out[209]: 4
```

Handling Missing Values

In [241]:

```
1  '''select missing values in a DataFrame. isnull and notnull return booleans
2  indicating whether a value is missing.'''
3  import pandas as pd
4  #Creating URL
5  url = 'C:/Users/Prerna/Desktop/ML_jupyter_notenooks/titanic.csv'
6
7  #Creating a dataframe
8  dataframe = pd.read_csv(url)
9
10 #Selecting the missing values from Age using isnull()
11
12 print(dataframe[dataframe['Age'].isnull()].head(5))
13
14 # Replacing Nan values to 0 using fillna()
15 #The fillna() is used to replace NaN values with zeros in Pandas DataFrame
16
17 dataframe = dataframe.fillna(0)
18
19 #checking is any null value left
20 print(dataframe.isnull().sum())
21
22 #checking is any null value left
23 print(dataframe.isnull().values.any())
24
25 #You can also replace Nan values with zeros when Loading the data
26 #dataframe = pd.read_csv(url).fillna(0)
27
```

	Name	PClass	Age	Sex	Survived	SexCode
12	Aubert, Mrs Leontine Pauline	1st	NaN	female	1	1
13	Barkworth, Mr Algernon H	1st	NaN	male	1	0
14	Baumann, Mr John D	1st	NaN	male	0	0
29	Borebank, Mr John James	1st	NaN	male	0	0
32	Bradley, Mr George	1st	NaN	male	1	0
	Name	0				
	PClass	0				
	Age	0				
	Sex	0				
	Survived	0				
	SexCode	0				

```
dtype: int64  
False
```

Dropping Duplicates

`drop_duplicates()` defaults to only dropping rows that match perfectly across all columns.

In [211]:

```
1 #Creating URL  
2 url = 'C:/Users/Prerna/Desktop/ML_jupyter_notenooks/titanic_duplicatevalues.csv'  
3  
4 #Creating a dataframe  
5 dataframe = pd.read_csv(url)  
6  
7 print("Original Data Frame:", len(dataframe))  
8 print("Dataframe after dropping duplicate values:",  
9       len(dataframe.drop_duplicates()))
```

Original Data Frame: 1323

Dataframe after dropping duplicate values: 1313

Grouping Rows by Values

`groupby` is one of the most powerful features in pandas.

`groupby` needs to be paired with some operation we want to apply to each group, such as calculating an aggregate statistic (e.g., mean, median, sum)

```
In [212]: 1 #Creating URL
          2 url = 'C:/Users/Prerna/Desktop/ML_jupyter_notenooks/titanic.csv'
          3
          4 #Creating a dataframe
          5 dataframe = pd.read_csv(url)
          6
          7 #Group by Gender
          8 dataframe.groupby('Sex').count()
```

Out[212]:

	Name	PClass	Age	Survived	SexCode
Sex					
<hr/>					
female	462	462	288	462	462
male	851	851	468	851	851

```
In [213]: 1 #Group by Gender
          2 dataframe.groupby('Sex')['Name'].count()
```

Out[213]: Sex
female 462
male 851
Name: Name, dtype: int64

In [214]:

```
1 #Group by multiple columns.  
2 '''Find the count of males and females by their survival wise and  
3 passenger class wise in titanic.'''  
4  
5 dataframe.groupby(['Sex', 'Survived', 'PClass'])['Name'].count()  
6
```

Out[214]:

Sex	Survived	PClass	
female	0	1st	9
		2nd	13
		3rd	132
	1	1st	134
		2nd	94
		3rd	80
male	0	*	1
		1st	120
		2nd	147
	1	3rd	441
		1st	59
		2nd	25
		3rd	58

Name: Name, dtype: int64

Looping Over a Column

To iterate over every element in a column and apply some action You can treat a pandas column like any other sequence in Python.

In [215]:

```
1 #Creating URL
2 url = 'C:/Users/Prerna/Desktop/ML_jupyter_notenooks/titanic.csv'
3
4 #Creating a dataframe
5 dataframe = pd.read_csv(url)
6
7 for name in dataframe['Name'][0:len(dataframe)]:
8     print(name.upper())
```

FARTHING, MR JOHN
FLEMING, MS MARGARET
FRANCATELLI, MS LAURA MABEL
FRY, MR RICHARD
GEIGER, MISS EMILY
GIGLIO, MR VICTOR
HARRINGTON, MR CHARLES
HARRISON, MR WILLIAM HENRY
HASSAH, MR HAMAD
ICABAD (ICABOD), MS
KEEPING, MR EDWIN
KENCHEN, MS AMELIA
LEROY, MISS BERTHE
LESNEUR, MR GUSTAVE
MALONEY, MS
OLIVA, MLLE
PERICAULT, MS
RINGHINI, MR SANTE
ROBBINS, MR VICTOR
SEGFESSER, MISS EMMA

Applying a Function Over All Elements in a Column

Use `apply()` to apply a built-in or custom function on every element in a column

In [216]:

```
1 #Creating URL
2 url = 'C:/Users/Prerna/Desktop/ML_jupyter_notenooks/titanic.csv'
3
4 #Creating a dataframe
5 dataframe = pd.read_csv(url)
6
7 #Create a userdefined function
8 def uppercase(x):
9     return x.upper()
10
11 #Apply custom function to the rows of the dataframe with apply()
12 dataframe['Name'].apply(uppercase)[0:len(dataframe)]
13
14
```

Out[216]:

```
0          ALLEN, MISS ELISABETH WALTON
1          ALLISON, MISS HELEN LORAIN
2          ALLISON, MR HUDSON JOSHUA CREIGHTON
3          ALLISON, MRS HUDSON JC (BESSIE WALDO DANIELS)
4          ALLISON, MASTER HUDSON TREVOR
5          ANDERSON, MR HARRY
6          ANDREWS, MISS KORNELIA THEODOSIA
7          ANDREWS, MR THOMAS, JR
8          APPLETON, MRS EDWARD DALE (CHARLOTTE LAMSON)
9          ARTAGAVEYTIA, MR RAMON
10         ASTOR, COLONEL JOHN JACOB
11         ASTOR, MRS JOHN JACOB (MADELEINE TALMADGE FORCE)
12         AUBERT, MRS LEONTINE PAULINE
13         BARKWORTH, MR ALGERNON H
14         BAUMANN, MR JOHN D
15         BAXTER, MRS JAMES (HELENE DELAUDENIERE CHAPUT)
16         BAXTER, MR QUIGG EDMOND
17         BEATTIE, MR THOMSON
18         BECKWITH, MR RICHARD LEONARD
19         BECKWITH, MRS RICHARD LEONARD (SALLIE MONYPENY)
20         BEHR, MR KARL HOWELL
21         BIRNBAUM, MR JAKOB
22         BISHOP, MR DICKINSON H
23         BISHOP, MRS DICKINSON H (HELEN WALTON)
24         BJORNSTRM-STEFFANSSON, MR MAURITZ HAKAN
```

25	BLACKWELL, MR STEPHEN WEART
26	BLANK, MR HENRY
27	BONNELL, MISS CAROLINE
28	BONNELL, MISS ELIZABETH
29	BOREBANK, MR JOHN JAMES
	...
1283	VESTROM, MISS HULDA AMANDA ADOLFINA
1284	VONK, MR JENKO
1285	WARE, MR FREDERICK
1286	WARREN, MR CHARLES WILLIAM
1287	WAZLI, MR YOUSIF
1288	WEBBER, MR JAMES
1289	WENNERSTROM, MR AUGUST EDVARD
1290	WENZEL, MR LINHART
1291	WIDEGREN, MR CHARLES PETER
1292	WIKLUND, MR JACOB ALFRED
1293	WILKES, MRS ELLEN
1294	WILLER, MR AARON
1295	WILLEY, MR EDWARD
1296	WILLIAMS, MR HOWARD HUGH
1297	WILLIAMS, MR LESLIE
1298	WINDELOV, MR EINAR
1299	WIRZ, MR ALBERT
1300	WISEMAN, MR PHILLIPPE
1301	WITTEVRONGEL, MR CAMIEL
1302	YALSEVAC, MR IVAN
1303	YASBECK, MR ANTONI
1304	YASBECK, MRS ANTONI
1305	YOUSSEF, MR GERIOS
1306	ZABOUR, MISS HILENI
1307	ZABOUR, MISS TAMINI
1308	ZAKARIAN, MR ARTUN
1309	ZAKARIAN, MR MAPRIEDER
1310	ZENNI, MR PHILIP
1311	LIEVENS, MR RENE
1312	ZIMMERMAN, LEO

Name: Name, Length: 1313, dtype: object

Applying a Function to Groups

You have grouped rows using groupby and want to apply a function to each group.

Combine groupby and apply.

In [217]:

```
1  #Creating URL
2  url = 'C:/Users/Prerna/Desktop/ML_jupyter_notenooks/titanic.csv'
3
4  #Creating a dataframe
5  dataframe = pd.read_csv(url)
6
7  # use groupby and apply
8
9  dataframe.groupby('Sex').apply(lambda x: x.count())
```

Out[217]:

	Name	PClass	Age	Sex	Survived	SexCode
Sex						
female	462	462	288	462	462	462
male	851	851	468	851	851	851

In [249]:

```
1  #DATA CLEANING CODE
2
3  import pandas as pd
4
5  #Creating URL
6  url = 'C:/Users/Prerna/Desktop/ML_jupyter_notenooks/titanic.csv'
7
8  #Creating a dataframe and replacing NAN values with 0
9  dataframe = pd.read_csv(url).fillna(0)
10
11
12  #checking is any null value left
13  print("Any null values in DataFrame:", dataframe.isnull().values.any())
14
15  #Drop duplicates.
16  dataframe.drop_duplicates()
17
18  #Drop column with drop() with axis=1 for column
19  dataframe = dataframe.drop('SexCode',axis=1)
20
21  #Rename column PClass to Passenger class using rename()
22  dataframe = dataframe.rename(columns = {'PClass': 'Passenger Class',
23                                         'Sex': 'Gender'})
24
25  #Replacing Survived column 0 and 1 with No and Yes
26  dataframe['Survived'] = dataframe['Survived'].replace([0,1],
27                                                         ['No', 'Yes'])
28
29  #Create a userdefined function
30  def uppercase(x):
31      return x.upper()
32
33  #Apply custom function to the rows of the dataframe with apply()
34  dataframe['Name'] = dataframe['Name'].apply(uppercase)[0:len(dataframe)]
35
36  dataframe
37
38  # Writing the dataframe to a csv file using to_csv(filename)
39  dataframe.to_csv('Cleareddata_titanic.csv')
40
```

Any null values in DataFrame: False

Out[249]:

	Name	Passenger Class	Age	Gender	Survived
0	ALLEN, MISS ELISABETH WALTON	1st	29.00	female	Yes
1	ALLISON, MISS HELEN LORAINE	1st	2.00	female	No
2	ALLISON, MR HUDSON JOSHUA CREIGHTON	1st	30.00	male	No
3	ALLISON, MRS HUDSON JC (BESSIE WALDO DANIELS)	1st	25.00	female	No
4	ALLISON, MASTER HUDSON TREVOR	1st	0.92	male	Yes
5	ANDERSON, MR HARRY	1st	47.00	male	Yes
6	ANDREWS, MISS KORNELIA THEODOSIA	1st	63.00	female	Yes
7	ANDREWS, MR THOMAS, JR	1st	39.00	male	No
8	APPLETON, MRS EDWARD DALE (CHARLOTTE LAMSON)	1st	58.00	female	Yes
9	ARTAGAVEYTIA, MR RAMON	1st	71.00	male	No
10	ASTOR, COLONEL JOHN JACOB	1st	47.00	male	No
11	ASTOR, MRS JOHN JACOB (MADELEINE TALMADGE FORCE)	1st	19.00	female	Yes
12	AUBERT, MRS LEONTINE PAULINE	1st	0.00	female	Yes
13	BARKWORTH, MR ALGERNON H	1st	0.00	male	Yes
14	BAUMANN, MR JOHN D	1st	0.00	male	No
15	BAXTER, MRS JAMES (HELENE DELAUDENIERE CHAPUT)	1st	50.00	female	Yes
16	BAXTER, MR QUIGG EDMOND	1st	24.00	male	No
17	BEATTIE, MR THOMSON	1st	36.00	male	No
18	BECKWITH, MR RICHARD LEONARD	1st	37.00	male	Yes
19	BECKWITH, MRS RICHARD LEONARD (SALLIE MONYPENY)	1st	47.00	female	Yes
20	BEHR, MR KARL HOWELL	1st	26.00	male	Yes
21	BIRNBAUM, MR JAKOB	1st	25.00	male	No
22	BISHOP, MR DICKINSON H	1st	25.00	male	Yes
23	BISHOP, MRS DICKINSON H (HELEN WALTON)	1st	19.00	female	Yes

	Name	Passenger Class	Age	Gender	Survived
24	BJORNSTRM-STEFFANSSON, MR MAURITZ HAKAN	1st	28.00	male	Yes
25	BLACKWELL, MR STEPHEN WEART	1st	45.00	male	No
26	BLANK, MR HENRY	1st	39.00	male	Yes
27	BONNELL, MISS CAROLINE	1st	30.00	female	Yes
28	BONNELL, MISS ELIZABETH	1st	58.00	female	Yes
29	BOREBANK, MR JOHN JAMES	1st	0.00	male	No
...
1283	VESTROM, MISS HULDA AMANDA ADOLFINA	3rd	14.00	female	No
1284	VONK, MR JENKO	3rd	22.00	male	No
1285	WARE, MR FREDERICK	3rd	0.00	male	No
1286	WARREN, MR CHARLES WILLIAM	3rd	0.00	male	No
1287	WAZLI, MR YOUSIF	3rd	0.00	male	No
1288	WEBBER, MR JAMES	3rd	0.00	male	No
1289	WENNERSTROM, MR AUGUST EDVARD	3rd	0.00	male	Yes
1290	WENZEL, MR LINHART	3rd	0.00	male	No
1291	WIDEGREN, MR CHARLES PETER	3rd	51.00	male	No
1292	WIKLUND, MR JACOB ALFRED	3rd	18.00	male	No
1293	WILKES, MRS ELLEN	3rd	45.00	female	Yes
1294	WILLER, MR AARON	3rd	0.00	male	No
1295	WILLEY, MR EDWARD	3rd	0.00	male	No
1296	WILLIAMS, MR HOWARD HUGH	3rd	0.00	male	No
1297	WILLIAMS, MR LESLIE	3rd	28.00	male	No
1298	WINDELOV, MR EINAR	3rd	21.00	male	No
1299	WIRZ, MR ALBERT	3rd	27.00	male	No
1300	WISEMAN, MR PHILLIPPE	3rd	0.00	male	No
1301	WITTEVRONGEL, MR CAMIEL	3rd	36.00	male	No

	Name	Passenger Class	Age	Gender	Survived
1302	YALSEVAC, MR IVAN	3rd	0.00	male	Yes
1303	YASBECK, MR ANTONI	3rd	27.00	male	No
1304	YASBECK, MRS ANTONI	3rd	15.00	female	Yes
1305	YOUSSEF, MR GERIOS	3rd	0.00	male	No
1306	ZABOUR, MISS HILENI	3rd	0.00	female	No
1307	ZABOUR, MISS TAMINI	3rd	0.00	female	No
1308	ZAKARIAN, MR ARTUN	3rd	27.00	male	No
1309	ZAKARIAN, MR MAPRIEDER	3rd	26.00	male	No
1310	ZENNI, MR PHILIP	3rd	22.00	male	No
1311	LIEVENS, MR RENE	3rd	24.00	male	No
1312	ZIMMERMAN, LEO	3rd	29.00	male	No

1312 rows x 6 columns

Concatenating DataFrames

You want to concatenate two DataFrames.

Use concat with axis=0 to concatenate along the row axis.

In []:

```
1  #creating url
2  url = 'C:/Users/Prerna/Desktop/ML_jupyter_notenooks/titanic.csv'
3
4  #Creating a dataframe and replacing NAN values with 0
5  dataframe = pd.read_csv(url).fillna(0)
6
7  # Creating new data frames
8  df1 = dataframe.iloc[0:5]
9  print(df1)
10 df2 = dataframe.iloc[5:10]
11 print(df2)
12
13 # Concetenating data Frames by row axis
14 df3 = pd.concat([df1,df2],axis=0)
15 print(df3)
16 # Concetenating data Frames by column axis
17 df4 = pd.concat([df1,df2],axis=1)
18 print(df4)
19
20
```

Merging DataFrames

To inner join, use merge with the on parameter to specify the column to merge on.

In [260]:

```
1 import pandas as pd
2
3 # Create DataFrame
4 employee_data = {'employee_id': ['1', '2', '3', '4'],
5 'name': ['Amy Jones', 'Allen Keys', 'Alice Bees',
6 'Tim Horton']}
7 dataframe_employees = pd.DataFrame(employee_data, columns = ['employee_id',
8 'name'])
9 # Create DataFrame
10 sales_data = {'employee_id': ['3', '4', '5', '6'],
11 'total_sales': [23456, 2512, 2345, 1455]}
12 dataframe_sales = pd.DataFrame(sales_data, columns = ['employee_id',
13 'total_sales'])
14 # Merge DataFrames using inner join
15 print(pd.merge(dataframe_employees, dataframe_sales, on='employee_id'))
16
17 # Merge DataFrames using outer join. Values for how = outer, left, right
18 print(pd.merge(dataframe_employees, dataframe_sales, on='employee_id', how='outer'))
19
20
```

	employee_id	name	total_sales
0	3	Alice Bees	23456
1	4	Tim Horton	2512

	employee_id	name	total_sales
0	1	Amy Jones	NaN
1	2	Allen Keys	NaN
2	3	Alice Bees	23456.0
3	4	Tim Horton	2512.0
4	5	NaN	2345.0
5	6	NaN	1455.0