```
 1 import pandas as pd
 2
 3 A Series is essentially a column, and a DataFrame is a
   multi-dimensional table made up of a collection of Series.
 4
 5 data = {
 6 'apples': [3, 2, 0, 1],
 7 'oranges': [0, 3, 7, 2]
 8 }
 9
10 purchases = pd.DataFrame(data)
11
12 purchases
13
14 purchases = pd.DataFrame(data, index=['June', 'Robert', '
   Lily', 'David'])
15
16 purchases
17
18 purchases.loc['June']
19
20 df = pd.read_csv('purchases.csv')
21
22 df
23
24 df = pd.read_csv('purchases.csv', index_col=0)
25
26 df
27
28 movies_df = pd.read_csv("IMDB-Movie-Data.csv", index_col="
   Title")
29
30 movies_df.head()
31
32 movies_df.tail(2)
33 movies_df.info()
34 movies_df.shape
35
36 temp_df = movies_df.append(movies_df)
37
38 temp_df.shape
39 temp_df.drop_duplicates(inplace=True)
40 temp_df = movies_df.append(movies_df) # make a new copy
41
42 temp_df.drop_duplicates(inplace=True, keep=False)
```

```python
43
44 temp_df.shape
45
46 #Column cleanup
47 #Many times datasets will have verbose column names with
   symbols, upper and lowercase words, spaces, and typos. To
   make selecting data by column name easier we can spend a
   little time cleaning up their names.
48
49 movies_df.columns
50 movies_df.rename(columns={
51 'Runtime (Minutes)': 'Runtime',
52 'Revenue (Millions)': 'Revenue_millions'
53 }, inplace=True)
54
55 movies_df.columns
56
57 movies_df.columns = ['rank', 'genre', 'description', '
   director', 'actors', 'year', 'runtime',
58 'rating', 'votes', 'revenue_millions', 'metascore']
59
60 movies_df.columns
61
62 movies_df.columns = [col.lower() for col in movies_df]
63
64 movies_df.columns
65
66 movies_df.isnull()
67 movies_df.isnull().sum()
68 movies_df.dropna()
69 movies_df.dropna(axis=1)
70 #What's with this axis=1parameter?
71
72 #It's not immediately obvious where axis comes from and
   why you need it to be 1 for it to affect columns. To see
   why, just look at the .shape output:
73
74 #movies_df.shape
75
76 #Out: (1000, 11)
77
78 #As we learned above, this is a tuple that represents the
   shape of the DataFrame, i.e. 1000 rows and 11 columns.
   Note that the rows are at index zero of this tuple and
   columns are at index one of this tuple. This is why axis=1
```

```python
78   affects columns. This comes from NumPy, and is a great
     example of why learning NumPy is worth your time.
79
80  #Imputation
81  #Imputation is a conventional feature engineering
    technique used to keep valuable data that have null
    values.
82
83  #Let's look at imputing the missing values in the
    revenue_millions column. First we'll extract that column
    into its own variable:
84
85  revenue = movies_df['revenue_millions']
86  revenue.head()
87  revenue_mean = revenue.mean()
88
89  revenue_mean
90  revenue.fillna(revenue_mean, inplace=True)
91  movies_df.isnull().sum()
92
93  #.describe() can also be used on a categorical variable
    to get the count of rows, unique count of categories, top
     category, and freq of top category:
94
95  movies_df.describe()
96  movies_df['genre'].describe()
97  #.value_counts() can tell us the frequency of all values
    in a column:
98
99  movies_df['genre'].value_counts().head(10)
100
101 #Relationships between continuous variables
102 #By using the correlation method .corr() we can generate
    the relationship between each continuous variable:
103
104 #movies_df.corr()
105
106 genre_col = movies_df['genre']
107
108 type(genre_col)
109
110 genre_col = movies_df[['genre']] #2
111
112 type(genre_col)
113
```

```
114 subset = movies_df[['genre', 'rating']]
115
116 subset.head()
117
118 For rows, we have two options:
119
120 #.loc - locates by name
121 #.iloc- locates by numerical index
122 #Remember that we are still indexed by movie Title, so to
      use .loc we give it the Title of a movie:
123
124 prom = movies_df.loc["Prometheus"]
125
126 prom
127
128 prom = movies_df.iloc[1]
129
130 movie_subset = movies_df.loc['Prometheus':'Sing']
131
132 movie_subset = movies_df.iloc[1:4]
133
134 movie_subset
135
136 condition = (movies_df['director'] == "Ridley Scott")
137
138 condition.head()
139
140 We want to filter out all movies not directed by Ridley
    Scott, in other words, we don't want the False films. To
    return the rows where that condition is True we have to
    pass this operation into the DataFrame:
141
142 movies_df[movies_df['director'] == "Ridley Scott"]
143 movies_df[movies_df['rating'] >= 8.6].head(3)
144 movies_df[(movies_df['director'] == 'Christopher Nolan')
    | (movies_df['director'] == 'Ridley Scott')].head()
145
146 def rating_function(x):
147 if x >= 8.0:
148 return "good"
149 else:
150 return "bad"
151
152 movies_df["rating_category"] = movies_df["rating"].apply(
    rating_function)
```

```
153
154 movies_df.head(2)
155
156 movies_df["rating_category"] = movies_df["rating"].apply(
    lambda x: 'good' if x >= 8.0 else 'bad')
157
158 movies_df.head(2)
159
160 import matplotlib.pyplot as plt
161 plt.rcParams.update({'font.size': 20, 'figure.figsize': (
    10, 8)}) # set font and plot size to be larger
162
163 movies_df.plot(kind='scatter', x='rating', y='
    revenue_millions', title='Revenue (millions) vs Rating');
164
165 movies_df['rating'].plot(kind="box");
166
```