

Data Wrangling

Agenda

In this session, we will cover the following concepts with the help of a business use case:

- Data acquisition
- Different methods for data wrangling:
 - Merge datasets
 - Concatenate datasets
 - Identify missing values
 - Drop unnecessary columns
 - Check the dimension of the dataset
 - Check the datatype of the dataset
 - Check datatype summary
 - Treat missing values
 - Validate correctness of the data in primary level if applicable

What Is Data Wrangling?

Data wrangling is the process of converting and formatting data from its raw form to usable format further down the data science pipeline.

What Is the Need for Data Wrangling?

Without feeding proper data into a model, one cannot expect a model that is dependable and gives higher accuracy.

Problem Statement

You are a junior data scientist and you are assigned a new task to perform data wrangling on a set of datasets. The datasets have many ambiguities. You have to identify those and apply different data wrangling techniques to get a dataset for further usage.

Dataset

- Download the `rental_bike_descr`, `rental_bike_season`, and `final_rental_bike_dataset` from Course Resources and upload the datasets to the lab

Data Dictionary

Attribute Information:

- date = date of the ride
- season - 1 = spring, 2 = summer, 3 = fall, 4 = winter
- holiday - whether the day is considered a holiday
- workingday - whether the day is neither a weekend nor holiday
- weather - 1: Clear, Few clouds, Partly cloudy, Partly cloudy
- 2: Mist + Cloudy, Mist + Broken clouds, Mist + Few clouds, Mist
- 3: Light Snow, Light Rain + Thunderstorm + Scattered clouds, Light Rain + Scattered clouds
- 4: Heavy Rain + Ice Pallets + Thunderstorm + Mist, Snow + Fog
- temp - temperature in Celsius
- atemp - "feels like" temperature in Celsius
- humidity - relative humidity
- windspeed - wind speed
- casual - number of non-registered user rentals initiated
- registered - number of registered user rentals initiated
- count - number of total rentals

Import libraries

- Pandas is a high-level data manipulation tool
- NumPy is used for working with multidimensional arrays

```
In [3]: import pandas as pd
import numpy as np
```

```
In [2]: print(pd.show_versions())

/usr/local/lib/python3.7/site-packages/pandas_datareader/compat/_init_.py:7: FutureWarning: pandas.util.testing is deprecated. Use the functions in the public API at pandas.testing instead.
  from pandas.util.testing import assert_frame_equal

INSTALLED VERSIONS
-----
commit           : b59586e1999e9aead1930c0bba2b674378607b3d
python           : 3.7.6.final.0
python-bits      : 64
OS              : Linux
OS-release       : 4.4.0-210-generic
Version         : #242~Ubuntu SMP Fri Apr 16 09:57:56 UTC 2021
processor        : x86_64
byteorder        : little
LC_ALL          : en_US.UTF-8
LANG            : en_US.UTF-8
LOCALE          : en_US.UTF-8

pandas          : 1.1.5
numpy           : 1.21.5
matplotlib      : 3.3.3
dateutil        : 2.0.9
pip            : 22.0.3
setuptools      : 43.0.0
Cython          : 0.29.16
pytest          : 5.4.1
hypothesis      : 5.9.0
ephem          : 2.4.4
klocwork        : 1.9.0
feather         : None
xlsxwriter      : 1.2.8
lxml.etree      : 4.5.0
traitlets       : 5.0.1
pymysql         : None
pymssql2        : None
cymysql         : 2.11.1
IPython         : 7.13.0
pandas_datareader: None
bs4             : 4.9.2
bottleneck      : 1.3.2
fsspec          : 0.8.3
fastparquet     : None
gcsfs           : None
matplotlib      : 3.3.1
numexpr         : 2.7.1
odfpy           : None
openpyxl        : 3.0.3
pandas_gbq      : None
pyarrow         : None
pytables        : None
pyxlsb          : None
s3fs            : None
scipy           : 1.4.1
sqlalchemy      : 1.3.15
tables          : 3.6.1
tabulate        : 0.8.7
xarray          : 0.15.1
xgboost         : 1.2.0
xlwt            : 1.3.0
numba          : 0.48.0
None
```

Load the first dataset

```
In [3]: dataset_1 = pd.read_csv('rental_bike_descr.csv')
```

Observations:

- We have to upload the dataset in the file explorer on the left panel of your lab
- We are reading the file through the `dataset_1` variable
- The file is in CSV format
- We use the `pd.read_csv()` function to read a CSV file
- We provide the exact path of the file within the round bracket ()

Check the type of dataset

- Execute the below command to understand type of data we are having

```
In [4]: type(dataset_1)

Out[4]: pandas.core.frame.DataFrame
```

Observations:

- The result shows that the dataset is DataFrame
- DataFrame is a tabular structure consisting of rows and columns

Shape of the dataset

```
In [5]: dataset_1.shape

Out[5]: (610, 10)
```

Observation:

- The `dataset_1` has 610 rows and 10 columns.

Print first 5 rows of the dataset

```
In [6]: dataset_1.head()
```

	instant	dteday	season	yr	mnth	hr	holiday	weekday	weathersit	temp
0	1	01-01-2011	1	0	1	0	False	6	1	0.24
1	2	01-01-2011	1	0	1	1	False	6	1	0.22
2	3	01-01-2011	1	0	1	2	False	6	1	0.22
3	4	01-01-2011	1	0	1	3	False	6	1	0.24
4	5	01-01-2011	1	0	1	4	False	6	1	0.24

Observation:

- The `'dataset_1.head()'` function displays only the initial five rows of the dataset.

Load the second dataset

- Use the function carefully since it is an excel file

```
In [7]: dataset_2 = pd.read_excel('rental_bike_season.xlsx')
```

Shape of the dataset

```
In [8]: dataset_2.shape

Out[8]: (610, 8)
```

Observation:

- The result shows that `dataset_2` has 610 rows and 8 columns.

Print first 5 rows of the dataset

```
In [9]: dataset_2.head()
```

	Unnamed: 0	instant	atemp	hum	windspeed	casual	registered	cnt
0	0	1	0.2879	0.81	0.0	3	13	16
1	1	2	0.2727	0.80	0.0	8	32	40
2	2	3	0.2727	0.80	0.0	5	27	32
3	3	4	0.2879	0.75	0.0	3	10	13
4	4	5	0.2879	0.75	0.0	0	1	1

Observation:

- We can see a column named `Unnamed: 0`, which is not in the data dictionary. Let's remove it.

Drop the column

```
In [10]: dataset_2 = dataset_2.drop(['Unnamed: 0'], axis=1)
```

Lets check the shape of the dataset again after the drop

```
In [11]: dataset_2.shape

Out[11]: (610, 7)
```

Observation:

- We had 8 columns before the drop.
- When we check the shape of the file after the drop, we see that the column `Unnamed: 0` has been dropped

Top 5 rows of the dataset

- Let's check the dataset_2 again

```
In [12]: dataset_2.head()
```

	instant	atemp	hum	windspeed	casual	registered	cnt
0	1	0.2879	0.81	0.0	3	13	16
1	2	0.2727	0.80	0.0	8	32	40
2	3	0.2727	0.80	0.0	5	27	32
3	4	0.2879	0.75	0.0	3	10	13
4	5	0.2879	0.75	0.0	0	1	1

Observation:

- dataset_2 does not have `Unnamed: 0` column

Merge the datasets

- We have two datasets. They are `dataset_1` and `dataset_2`
- As both datasets have one common column `instant`, let's merge the datasets on that column
- We are going to save the resultant data inside the `combined_data` as shown below

```
In [13]: combined_data = pd.merge(dataset_1, dataset_2, on='instant')
```

Check the shape of combined dataset

```
In [14]: combined_data.shape

Out[14]: (610, 16)
```

Observation:

- The shape of the combined_data has 610 rows and 16 columns

Top 5 rows of the combined dataset

```
In [15]: combined_data.head()
```

	instant	dteday	season	yr	mnth	hr	holiday	weekday	weathersit	temp	atemp	hum	windspeed	casual	registered	cnt
0	1	01-01-2011	1	0	1	0	False	6	1	0.24	0.2879	0.81	0.0	3	13	16
1	2	01-01-2011	1	0	1	1	False	6	1	0.22	0.2727	0.80	0.0	8	32	40
2	3	01-01-2011	1	0	1	2	False	6	1	0.22	0.2727	0.80	0.0	5	27	32
3	4	01-01-2011	1	0	1	3	False	6	1	0.24	0.2879	0.75	0.0	3	10	13
4	5	01-01-2011	1	0	1	4	False	6	1	0.24	0.2879	0.75	0.0	0	1	1

Now, load the 3rd dataset

- The dataset is saved in s3 bucket, we are going to download the dataset_3

Load the third dataset

Import the dataset

```
In [18]: dataset_3 = pd.read_csv('final_rental_bike_dataset.csv')
```

Check the shape of the dataset

```
In [19]: dataset_3.shape

Out[19]: (390, 16)
```

Top 5 rows of the dataset

```
In [20]: dataset_3.head()
```

	instant	dteday	season	yr	mnth	hr	holiday	weekday	weathersit	temp	atemp	hum	windspeed	casual	registered	cnt
0	620	29-01-2011	1	0	1	0	False	6	1	0.36	0.3333	0.40	0.2985	0	2	2
1	621	29-01-2011	1	0	1	1	False	6	1	0.34	0.3182	0.64	0.1940	0	10	15
2	622	29-01-2011	1	0	1	2	False	6	1	0.20	0.2121	0.64	0.1343	3	5	8
3	623	29-01-2011	1	0	1	3	False	6	1	0.16	0.1818	0.69	0.1045	1	2	3
4	624	29-01-2011	1	0	1	4	False	6	1	0.16	0.1818	0.64	0.1343	0	2	2

Bottom 15 rows of the dataset

- Just like the `head` function, the `tail` function is used to see the bottom rows of the dataset
- If you want to see the specific number of rows, then specify the number inside the `bracket ()` as shown below

```
In [21]: dataset_3.tail(15)
```

	instant	dteday	season	yr	mnth	hr	holiday	weekday	weathersit	temp	atemp	hum	windspeed	casual	registered	cnt
375	995	14-02-2011	1	0	2	2	False	1	1	0.36	0.3333	0.40	0.2985	0	2	2
376	996	14-02-2011	1	0	2	3	False	1	1	0.34	0.3182	0.64	0.1940	1	1	2
377	997	14-02-2011	1	0	2	4	False	1	1	0.32	0.3030	0.53	0.2836	0	2	2
378	998	14-02-2011	1	0	2	5	False	1	1	0.32	0.3030	0.53	0.2836	0	3	3
379	999	14-02-2011	1	0	2	6	False	1	1	0.34	0.3030	0.46	0.2985	1	25	26
380	1000	14-02-2011	1	0	2	7	False	1	1	0.34	0.3030	0.46	0.2985	2	96	98
381	611	28-01-2011	1	0	1	16	False	5	1	0.22	0.2727	0.80	0.0000	10	70	80
382	612	28-01-2011	1	0	1	17	False	5	1	0.24	0.2424	0.75	0.1343	2	147	149
383	613	28-01-2011	1	0	1	18	False	5	1	0.24	0.2273	0.75	0.1940	2	107	109
384	614	28-01-2011	1	0	1	19	False	5	2	0.24	0.2424	0.75	0.1343	5	84	89
385	615	28-01-2011	1	0	1	20	False	5	2	0.24	0.2273	0.70	0.1940	1	61	62
386	616	28-01-2011	1	0	1	21	False	5	2	0.22	0.2273	0.75	0.1343	1	57	58
387	617	28-01-2011	1	0	1	22	False	5	1	0.24	0.2121	0.65	0.3582	0	26	26
388	618	28-01-2011	1	0	1	23	False	5	1	0.24	0.2273	0.60	0.2239	1	22	23
389	619	29-01-2011	1	0	1	0	False	6	1	0.22	0.1970	0.64	0.3582	2	26	28

Observation:

- The bottom 15 rows of the dataset_3 is shown above, as we mention 15 inside the bracket ()
- Here, we can see that the rows are not sorted well according to the `instant` number. Let's resolve it.

Sort values of a column

- To sort the values per our will, we use the `sort_values` function and in the square brackets, we specify the name of the column by which we want to sort it, as shown below

```
In [22]: dataset_3 = dataset_3.sort_values(by=['instant'])
```

- Let's check head and tail to verify the sort operation

```
In [23]: dataset_3.head()
```

	instant	dteday	season	yr	mnth	hr	holiday	weekday	weathersit	temp	atemp	hum	windspeed	casual	registered	cnt
381	611	28-01-2011	1	0	1	16	False	5	1	0.22	0.2727	0.80	0.0000	10	70	80
382	612	28-01-2011	1	0	1	17	False	5	1	0.24	0.2424	0.75	0.1343	2	147	149
383	613	28-01-2011	1	0	1	18	False	5	1	0.24	0.2273	0.75	0.1940	2	107	109
384	614	28-01-2011	1	0	1	19	False	5	2	0.24	0.2424	0.75	0.1343	5	84	89
385	615	28-01-2011	1	0	1	20	False	5	2	0.24	0.2273	0.70	0.1940	1	61	62

```
In [24]: dataset_3.tail()
```

	instant	dteday	season	yr	mnth	hr	holiday	weekday	weathersit	temp	atemp	hum	windspeed	casual	registered	cnt
376	996	14-02-2011	1	0	2	3	False	1	1	0.34	0.3182	0.64	0.2239	1	1	2
377	997	14-02-2011	1	0	2	4	False	1	1	0.32	0.3030	0.53	0.2836	0	2	2
378	998	14-02-2011	1	0	2	5	False	1	1	0.34	0.3030	0.53	0.2836	0	3	3
379	999	14-02-2011	1	0	2	6	False	1	1	0.34	0.3030	0.46	0.2985	1	25	26
380	1000	14-02-2011	1	0	2	7	False	1	1	0.34	0.3030	0.46	0.2985	2	96	98

Concatenate the combine_data with dataset_3

- Let's concatenate both DataFrame `combined_data` and `dataset_3` into a single DataFrame using the `concat` function, as shown below
- Store the final DataFrame inside the `final_data` variable


```
... Boston dataset:
Boston house prices dataset

**Data Set Characteristics:**
 :Number of Instances: 506
 :Number of Attributes: 13 numeric/categorical predictive. Median Value (attribute 14) is usually the target variable.
 :Attribute Information (in order):
  - CRIM      per capita crime rate by town
  - ZN        proportion of residential land zoned for lots over 25,000 sq.ft.
  - INDUS     proportion of non-retail business acres per town
  - CHAS      Charles River dummy variable (= 1 if tract bounds river; 0 otherwise)
  - NOX       nitric oxides concentration (parts per 10 million)
  - RM        average number of rooms per dwelling
  - AGE       proportion of owner-occupied units built prior to 1940
  - DIS       weighted distances to five Boston employment centres
  - RAD       index of accessibility to radial highways
  - TAX       full-value property-tax rate per $10,000
  - PTRATIO   pupil-teacher ratio by town
  - B         1000(BK - 0.63)^2 where BK is the proportion of black people by town
  - LSTAT     lower status of the population
  - MEDV      Median value of owner-occupied homes in $1000's

 :Missing Attribute Values: None

 :Creator: Harrison, D. and Rubinfeld, D.L.

This is a copy of UCI ML housing dataset.
https://archive.ics.uci.edu/ml/machine-learning-databases/housing/

The dataset was taken from the StatLib library which is maintained at Carnegie Mellon University.

The Boston house-price data of Harrison, D. and Rubinfeld, D.L. "Hedonic
prices and the demand for clean air", J. Environ. Economics & Management,
vol.5, 41-102, 1981. Used in Belisle, R. Kuh & Welch, "Regression diagnostics
...", Wiley, 1980. N.B. Various transformations are used in the table on
pages 244-261 of the latter.

The Boston house-price data has been used in many machine learning papers that address regression
problems.

... topic1: References
  - Belisle, Kuh & Welch, "Regression diagnostics: Identifying Influential Data and Sources of Collinearity",
    Wiley, 1980. 244-261.
  - Quinlan, R. (1993). Combining Instance-Based and Model-Based Learning. In Proceedings on the Tenth Internat
    tional Conference of Machine Learning, 236-243, University of Massachusetts, Amherst. Morgan Kaufmann.
```

Let's convert the array to a DataFrame i.e into tabular structure

- Since both data and target are in a NumPy array, we need to convert it to a DataFrame

```
In [43]: boston_df = pd.DataFrame(boston_data)
boston_df.columns = columns
boston_df['MEDV'] = y
boston_df.o = boston_df
print(boston_df.shape)
print(boston_df.o.shape)

(506, 14)
(506, 14)
```

Observation:

- boston_df has 506 rows and 14 columns

```
In [44]: print(boston_df.describe())

count    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000
mean      3.613524    11.363636    11.136779    0.069170    0.556951    0.329664    -0.556961    0.632905    0.307898
std       8.601545    23.322453    6.860353    0.253994    0.115878    0.115878    0.115878    0.115878
min       0.006320    0.000000    0.460000    0.000000    0.385000    0.385000    0.385000    0.385000
25%       0.082045    0.000000    0.519000    0.000000    0.449000    0.449000    0.449000    0.449000
50%       0.256150    0.000000    0.690000    0.000000    0.538000    0.538000    0.538000    0.538000
75%       3.677083    12.500000    18.100000    0.000000    0.624000    0.624000    0.624000    0.624000
max       88.976200    100.000000    27.740000    1.000000    0.871000    0.871000    0.871000    0.871000
```

```
In [45]: boston_df.head()

count    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000
mean      68.574901    3.795043    9.549407    408.237154    18.455334    356.674032    18.455334    356.674032
std      28.148861    2.105710    1.807259    168.537116    2.164946    91.294864    2.164946    91.294864
min       2.900000    1.129600    1.000000    187.000000    12.600000    0.320000    12.600000    0.320000
25%      45.025000    1.299750    1.000000    279.000000    17.400000    375.377500    17.400000    375.377500
50%      77.500000    3.207450    5.000000    330.000000    19.050000    391.440000    19.050000    391.440000
75%      94.075000    5.188425    24.000000    666.000000    20.200000    396.225000    20.200000    396.225000
max     100.000000    12.125000    24.000000    711.000000    22.000000    396.900000    22.000000    396.900000
```

```
In [46]: boston_df.head()

count    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000
mean      68.574901    3.795043    9.549407    408.237154    18.455334    356.674032    18.455334    356.674032
std      28.148861    2.105710    1.807259    168.537116    2.164946    91.294864    2.164946    91.294864
min       2.900000    1.129600    1.000000    187.000000    12.600000    0.320000    12.600000    0.320000
25%      45.025000    1.299750    1.000000    279.000000    17.400000    375.377500    17.400000    375.377500
50%      77.500000    3.207450    5.000000    330.000000    19.050000    391.440000    19.050000    391.440000
75%      94.075000    5.188425    24.000000    666.000000    20.200000    396.225000    20.200000    396.225000
max     100.000000    12.125000    24.000000    711.000000    22.000000    396.900000    22.000000    396.900000
```

```
In [47]: boston_df.head()

count    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000
mean      68.574901    3.795043    9.549407    408.237154    18.455334    356.674032    18.455334    356.674032
std      28.148861    2.105710    1.807259    168.537116    2.164946    91.294864    2.164946    91.294864
min       2.900000    1.129600    1.000000    187.000000    12.600000    0.320000    12.600000    0.320000
25%      45.025000    1.299750    1.000000    279.000000    17.400000    375.377500    17.400000    375.377500
50%      77.500000    3.207450    5.000000    330.000000    19.050000    391.440000    19.050000    391.440000
75%      94.075000    5.188425    24.000000    666.000000    20.200000    396.225000    20.200000    396.225000
max     100.000000    12.125000    24.000000    711.000000    22.000000    396.900000    22.000000    396.900000
```

```
In [48]: boston_df.head()

count    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000
mean      68.574901    3.795043    9.549407    408.237154    18.455334    356.674032    18.455334    356.674032
std      28.148861    2.105710    1.807259    168.537116    2.164946    91.294864    2.164946    91.294864
min       2.900000    1.129600    1.000000    187.000000    12.600000    0.320000    12.600000    0.320000
25%      45.025000    1.299750    1.000000    279.000000    17.400000    375.377500    17.400000    375.377500
50%      77.500000    3.207450    5.000000    330.000000    19.050000    391.440000    19.050000    391.440000
75%      94.075000    5.188425    24.000000    666.000000    20.200000    396.225000    20.200000    396.225000
max     100.000000    12.125000    24.000000    711.000000    22.000000    396.900000    22.000000    396.900000
```

```
In [49]: boston_df.head()

count    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000
mean      68.574901    3.795043    9.549407    408.237154    18.455334    356.674032    18.455334    356.674032
std      28.148861    2.105710    1.807259    168.537116    2.164946    91.294864    2.164946    91.294864
min       2.900000    1.129600    1.000000    187.000000    12.600000    0.320000    12.600000    0.320000
25%      45.025000    1.299750    1.000000    279.000000    17.400000    375.377500    17.400000    375.377500
50%      77.500000    3.207450    5.000000    330.000000    19.050000    391.440000    19.050000    391.440000
75%      94.075000    5.188425    24.000000    666.000000    20.200000    396.225000    20.200000    396.225000
max     100.000000    12.125000    24.000000    711.000000    22.000000    396.900000    22.000000    396.900000
```

```
In [50]: boston_df.head()

count    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000
mean      68.574901    3.795043    9.549407    408.237154    18.455334    356.674032    18.455334    356.674032
std      28.148861    2.105710    1.807259    168.537116    2.164946    91.294864    2.164946    91.294864
min       2.900000    1.129600    1.000000    187.000000    12.600000    0.320000    12.600000    0.320000
25%      45.025000    1.299750    1.000000    279.000000    17.400000    375.377500    17.400000    375.377500
50%      77.500000    3.207450    5.000000    330.000000    19.050000    391.440000    19.050000    391.440000
75%      94.075000    5.188425    24.000000    666.000000    20.200000    396.225000    20.200000    396.225000
max     100.000000    12.125000    24.000000    711.000000    22.000000    396.900000    22.000000    396.900000
```

```
In [51]: boston_df.head()

count    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000
mean      68.574901    3.795043    9.549407    408.237154    18.455334    356.674032    18.455334    356.674032
std      28.148861    2.105710    1.807259    168.537116    2.164946    91.294864    2.164946    91.294864
min       2.900000    1.129600    1.000000    187.000000    12.600000    0.320000    12.600000    0.320000
25%      45.025000    1.299750    1.000000    279.000000    17.400000    375.377500    17.400000    375.377500
50%      77.500000    3.207450    5.000000    330.000000    19.050000    391.440000    19.050000    391.440000
75%      94.075000    5.188425    24.000000    666.000000    20.200000    396.225000    20.200000    396.225000
max     100.000000    12.125000    24.000000    711.000000    22.000000    396.900000    22.000000    396.900000
```

```
In [52]: boston_df.head()

count    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000
mean      68.574901    3.795043    9.549407    408.237154    18.455334    356.674032    18.455334    356.674032
std      28.148861    2.105710    1.807259    168.537116    2.164946    91.294864    2.164946    91.294864
min       2.900000    1.129600    1.000000    187.000000    12.600000    0.320000    12.600000    0.320000
25%      45.025000    1.299750    1.000000    279.000000    17.400000    375.377500    17.400000    375.377500
50%      77.500000    3.207450    5.000000    330.000000    19.050000    391.440000    19.050000    391.440000
75%      94.075000    5.188425    24.000000    666.000000    20.200000    396.225000    20.200000    396.225000
max     100.000000    12.125000    24.000000    711.000000    22.000000    396.900000    22.000000    396.900000
```

```
In [53]: boston_df.head()

count    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000
mean      68.574901    3.795043    9.549407    408.237154    18.455334    356.674032    18.455334    356.674032
std      28.148861    2.105710    1.807259    168.537116    2.164946    91.294864    2.164946    91.294864
min       2.900000    1.129600    1.000000    187.000000    12.600000    0.320000    12.600000    0.320000
25%      45.025000    1.299750    1.000000    279.000000    17.400000    375.377500    17.400000    375.377500
50%      77.500000    3.207450    5.000000    330.000000    19.050000    391.440000    19.050000    391.440000
75%      94.075000    5.188425    24.000000    666.000000    20.200000    396.225000    20.200000    396.225000
max     100.000000    12.125000    24.000000    711.000000    22.000000    396.900000    22.000000    396.900000
```

```
In [54]: boston_df.head()

count    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000
mean      68.574901    3.795043    9.549407    408.237154    18.455334    356.674032    18.455334    356.674032
std      28.148861    2.105710    1.807259    168.537116    2.164946    91.294864    2.164946    91.294864
min       2.900000    1.129600    1.000000    187.000000    12.600000    0.320000    12.600000    0.320000
25%      45.025000    1.299750    1.000000    279.000000    17.400000    375.377500    17.400000    375.377500
50%      77.500000    3.207450    5.000000    330.000000    19.050000    391.440000    19.050000    391.440000
75%      94.075000    5.188425    24.000000    666.000000    20.200000    396.225000    20.200000    396.225000
max     100.000000    12.125000    24.000000    711.000000    22.000000    396.900000    22.000000    396.900000
```

```
In [55]: boston_df.head()

count    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000
mean      68.574901    3.795043    9.549407    408.237154    18.455334    356.674032    18.455334    356.674032
std      28.148861    2.105710    1.807259    168.537116    2.164946    91.294864    2.164946    91.294864
min       2.900000    1.129600    1.000000    187.000000    12.600000    0.320000    12.600000    0.320000
25%      45.025000    1.299750    1.000000    279.000000    17.400000    375.377500    17.400000    375.377500
50%      77.500000    3.207450    5.000000    330.000000    19.050000    391.440000    19.050000    391.440000
75%      94.075000    5.188425    24.000000    666.000000    20.200000    396.225000    20.200000    396.225000
max     100.000000    12.125000    24.000000    711.000000    22.000000    396.900000    22.000000    396.900000
```

```
In [56]: boston_df.head()

count    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000
mean      68.574901    3.795043    9.549407    408.237154    18.455334    356.674032    18.455334    356.674032
std      28.148861    2.105710    1.807259    168.537116    2.164946    91.294864    2.164946    91.294864
min       2.900000    1.129600    1.000000    187.000000    12.600000    0.320000    12.600000    0.320000
25%      45.025000    1.299750    1.000000    279.000000    17.400000    375.377500    17.400000    375.377500
50%      77.500000    3.207450    5.000000    330.000000    19.050000    391.440000    19.050000    391.440000
75%      94.075000    5.188425    24.000000    666.000000    20.200000    396.225000    20.200000    396.225000
max     100.000000    12.125000    24.000000    711.000000    22.000000    396.900000    22.000000    396.900000
```

```
In [57]: boston_df.head()

count    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000
mean      68.574901    3.795043    9.549407    408.237154    18.455334    356.674032    18.455334    356.674032
std      28.148861    2.105710    1.807259    168.537116    2.164946    91.294864    2.164946    91.294864
min       2.900000    1.129600    1.000000    187.000000    12.600000    0.320000    12.600000    0.320000
25%      45.025000    1.299750    1.000000    279.000000    17.400000    375.377500    17.400000    375.377500
50%      77.500000    3.207450    5.000000    330.000000    19.050000    391.440000    19.050000    391.440000
75%      94.075000    5.188425    24.000000    666.000000    20.200000    396.225000    20.200000    396.225000
max     100.000000    12.125000    24.000000    711.000000    22.000000    396.900000    22.000000    396.900000
```

```
In [58]: boston_df.head()

count    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000    506.000000
mean      68.574901    3.795043    9.549407    408.237154    18.455334    356.674032    18.455334    356.674032
std      28.148861    2.105710    1.807259    168.537116    2.164946    91.294864    2.164946    91.294864
min       2.900000    1.129600    1.000000    187.000000    12.600000    0.320000    12.600000    0.320000
25%      45.025000    1.299750    1.000000    279.000000    17.400000    375.377500    17.400000    375.377500
50%      77.500000    3.207450    5.000000    330.000000    19.050000    391.440000    19.050000    391.440000
75%      94.075000    5.188425    24.000000    666.000000    20.200000    396.225000    20.200000    396.225000
max     100.00
```