

```
import pandas as pd
train = pd.read_csv("train-data.csv")

train.info() # info() 함수는 데이터에 대한 전반적인 정보를 나타냅니다.
#특성 구성에는 행과 열의 크기, 결측값, 데이터 타입을 구성하는 값의 자료형 등을 출력해줍니다.

<class 'pandas.core.frame.DataFrame'>
Int64Index: 6019 entries, 0 to 6018
Data columns (total 14 columns):
 #   Column                Non-Null Count  Dtype  ---  --
 0   Name                   6019 non-null    object  1   Name                   6019 non-null    object  2   Location               6019 non-null    object  3   Year                   6019 non-null    int64   4   Fuel_Type              6019 non-null    object  5   Transmission            6019 non-null    object  6   Mileage                 6017 non-null    object  7   Engine                 6019 non-null    object  8   Power                  6017 non-null    object  9   Seats                  6019 non-null    int64   10  New_Price              624 non-null     object  11  Price                 6019 non-null    float64 12  Price                 824 non-null     object 13  Price                 6019 non-null    float64
dtypes: float64(2), int64(3), object(9)
memory usage: 608.5+ KB

데이터 분석

14개의 항목, 6019개의 데이터, 행인 마력, 신차가격, 좌석 수에 걸쳐있다.

8 Mileage 6017 non-null object
9년 Mileage 은 2328로 부족하다.

9 Engine 5983 non-null object

10 Power 5983 non-null object

11 Seats 5977 non-null float64

12 New_Price 824 non-null object
신차 가격은 결측치가 많다.

13 Price 6019 non-null float64

Name의 특성부터 치러, Name은 행이만만 남기고 뒤의 글자는 삭제한다.
```

Owner_Type 전처리 & csv 파일에서 인덱스 삭제

```
In [2]: train['Owner_Type'] = train['Owner_Type'].replace('First', 1)
train['Owner_Type'] = train['Owner_Type'].replace('Second', 2)
train['Owner_Type'] = train['Owner_Type'].replace('Third', 3)
train['Owner_Type'] = train['Owner_Type'].replace('Fourth & Above', 4)
train = train.drop('Unnamed: 0', axis=1)
train.replace를 하면, replace를 쓰면서 다시 옮겨와야한다. 그렇지 않으면 데이터 양이 바뀌게 된다.
#데이터 자체에 인덱스가 있어서 해당 값을 지워야한다.

Out[2]:
```

	Name	Location	Year	Kilometers_Driven	Fuel_Type	Transmission	Owner_Type	Mileage	Engine	Power	Seats	New_Price	Price
0	Maruti Wagon R Lxi CNG	Mumbai	2010	72000	CNG	Manual	1	26.6 kmpl	998 CC	58.16 bhp	5.0	NaN	1.75
1	Hyundai Creta LE CRDI SX Opton	Pune	2015	41000	Diesel	Manual	1	19.67 kmpl	1552 CC	126.2 bhp	5.0	NaN	12.50
2	Honda Jazz V	Chennai	2011	46000	Petrol	Manual	1	18.2 kmpl	1196 CC	88.7 bhp	5.0	8.61 Lakh	4.50
3	Maruti Ertiga VDI	Chennai	2012	87000	Diesel	Manual	1	20.77 kmpl	1248 CC	88.76 bhp	7.0	NaN	6.00
4	Audi A4 New 2.0 T Multitronic	Coimbatore	2012	40670	Diesel	Automatic	2	15.2 kmpl	1968 CC	140.8 bhp	5.0	NaN	17.44
...
6014	Maruti Swift VDI	Delhi	2014	27365	Diesel	Manual	1	28.4 kmpl	1248 CC	74 bhp	5.0	7.88 Lakh	4.75
6015	Hyundai Xcent i11 CRDI S	Japur	2015	100000	Diesel	Manual	1	24.4 kmpl	1120 CC	71 bhp	5.0	NaN	2.90
6016	Mahindra Xeno i4 BSIV	Japur	2012	55000	Diesel	Manual	2	14.0 kmpl	2498 CC	112 bhp	8.0	NaN	2.90
6017	Maruti Wagon R VDI	Kolkata	2013	46000	Petrol	Manual	1	0.0 kmpl	2112 CC	nan bhp	6.0	NaN	2.65
6018	Chevrolet Beat Diesel	Hyderabad	2013	47000	Diesel	Manual	1	25.44 kmpl	936 CC	57.6 bhp	5.0	NaN	2.50

6019 rows x 13 columns

Name 전처리

```
In [3]: train['Name'] = train['Name'].str.extract('([A-Z]+)')
#([A-Z]+)는 문자를 제외한 공백없이 한글자를 추출, <+가 붙어서 단어 단위로 추출한다.
train['Name'] = train['Name'].str.extract('([A-Z]+)')

Mileage 전처리

In [4]: train[train['Mileage'].isnull() == True]
#train에 마일리지 값이 없는 데이터가 남아있다.

Out[4]:
```

	Name	Location	Year	Kilometers_Driven	Fuel_Type	Transmission	Owner_Type	Mileage	Engine	Power	Seats	New_Price	Price
4446	Mahindra	Chennai	2016	50000	Electric	Automatic	1	NaN	72 CC	41 bhp	5.0	13.58 Lakh	13.00
4904	Toyota	Mumbai	2011	44000	Electric	Automatic	1	NaN	1798 CC	73 bhp	5.0	NaN	32.75

```
In [5]: train[train['Fuel_Type'] == "Electric"]

Out[5]:
```

	Name	Location	Year	Kilometers_Driven	Fuel_Type	Transmission	Owner_Type	Mileage	Engine	Power	Seats	New_Price	Price
4446	Mahindra	Chennai	2016	50000	Electric	Automatic	1	NaN	72 CC	41 bhp	5.0	13.58 Lakh	13.00
4904	Toyota	Mumbai	2011	44000	Electric	Automatic	1	NaN	1798 CC	73 bhp	5.0	NaN	32.75

Mahindra E verito D4 판매, 전기차이다.

따라서 연비가 리터당 마일 아니라 거리로 때문에 연비가 잘못으로 나타난 것이다.

이런 경우 두 가지가 한 개뿐이어서 전체 모델에 미치는 영향은 아주 미미할 것으로 생각했다. 따라서 마힌다의 차량의 데이터는 지우기로 결정했다.

반면, 토요타 프리우스 2011년식은 29.3으로 연비가 존재한다. 따라서 이 경우에는 연비를 새로 입력했다.

```
In [6]: train.loc[4904, 'Mileage'] == 29.3
#마일리지 29.3은 토요타 프리우스 2011년식의 연비이다.
train['Mileage'].dropna()
#마일리지 값이 없는 데이터는 삭제한다.

Out[6]:
```

	Name	Location	Year	Kilometers_Driven	Fuel_Type	Transmission	Owner_Type	Mileage	Engine	Power	Seats	New_Price	Price
0	26.6 kmpl												
1	19.67 kmpl												
2	18.2 kmpl												
3	20.77 kmpl												
4	15.2 kmpl												
6014	28.4 kmpl												
6015	24.4 kmpl												
6016	14.0 kmpl												
6017	0.0 kmpl												
6018	25.44 kmpl												
6019	29.3 kmpl												

6019 rows x 13 columns

```
In [7]: train[train['Mileage'] == 0.0 kmpl]
#마일리지 0.0인 데이터는 마힌다의 차량들이 있다.
#마일리지 이상치들을 우선 연비의 중앙값으로 채운 뒤, 나중에 모델의 성능이 떨어지면
#그때 이 데이터를 삭제한다. (이 과정에서 정확도가 떨어질 것 같다.)

Out[7]:
```

	Name	Location	Year	Kilometers_Driven	Fuel_Type	Transmission	Owner_Type	Mileage	Engine	Power	Seats	New_Price	Price
14	Land	Pune	2012	85000	Diesel	Automatic	2	0.0 kmpl	2179 CC	113 bhp	5.0	NaN	17.50
67	Mercedes	Coimbatore	2010	15080	Diesel	Automatic	1	0.0 kmpl	1950 CC	194 bhp	5.0	49.14 Lakh	36.67
79	Hyundai	Hyderabad	2005	87561	Petrol	Manual	1	0.0 kmpl	1086 CC	nan bhp	5.0	NaN	1.30
134	Hyda	Ahmedabad	2007	60006	Petrol	Manual	1	0.0 kmpl	NaN	NaN	NaN	NaN	2.95
229	Ford	Bangalore	2015	70436	Diesel	Manual	1	0.0 kmpl	1491 CC	99 bhp	NaN	NaN	3.60
...
5647	Toyota	Mumbai	2011	227000	Diesel	Manual	4	0.0 kmpl	2444 CC	nan bhp	8.0	NaN	2.20
5875	Mercedes	Ahmedabad	2019	4000	Diesel	Automatic	1	0.0 kmpl	1950 CC	194 bhp	5.0	49.14 Lakh	35.00
5943	Mahindra	Chennai	2002	75000	Diesel	Manual	1	0.0 kmpl	2112 CC	nan bhp	6.0	NaN	17.00
5972	Hyundai	Mumbai	2008	65000	Petrol	Manual	1	0.0 kmpl	1086 CC	62 bhp	5.0	NaN	1.39
6011	Skoda	Hyderabad	2009	53000	Petrol	Automatic	1	0.0 kmpl	3597 CC	262.6 bhp	5.0	NaN	4.75

60 rows x 13 columns

오차 발생 가능성

아래 출력에서 보듯이, 연비가 0.0kmpl인 차량들이 있다.

이러한 이상치들은 우선 연비의 중앙값으로 채운 뒤, 나중에 모델의 성능이 떨어지면

그때 더 자세한 전처리를 해보도록 할 것이다.

```
In [8]: import numpy as np
train['Mileage'] = train['Mileage'].str.extract('([0-9]+\.?[0-9]+)')
train['Mileage'] = train['Mileage'].astype(float)
train['Mileage'] = train['Mileage'].replace(0, train['Mileage'].median())

#연비 한계치 '0.8kmpl' 이상인 것들을, 우선 float형으로 바꾸고,
#연비의 중앙값을 대입한다. (이 과정에서 정확도가 떨어질 것 같다.)

Engine 전처리

In [9]: train[train['Engine'].isnull() == True]
#마일리지 0.0인 데이터는 마힌다의 차량들이 있다.
#연비 한계치 '0.8kmpl' 이상인 것들을, 우선 float형으로 바꾸고,
#연비의 중앙값을 대입한다. (이 과정에서 정확도가 떨어질 것 같다.)

Out[9]:
```

	Name	Location	Year	Kilometers_Driven	Fuel_Type	Transmission	Owner_Type	Mileage	Engine	Power	Seats	New_Price	Price
194	Honda	Ahmedabad	2007	60006	Petrol	Manual	1	18.15	NaN	NaN	NaN	NaN	2.95
208	Maruti	Kolkata	2010	42001	Petrol	Manual	1	16.10	NaN	NaN	NaN	NaN	2.11
733	Maruti	Chennai	2006	97890	Petrol	Manual	3	16.10	NaN	NaN	NaN	NaN	1.75
749	Land	Mumbai	2008	55001	Diesel	Automatic	2	18.15	NaN	NaN	NaN	NaN	26.50
1294	Honda	Delhi	2009	55005	Petrol	Manual	1	12.90	NaN	NaN	NaN	NaN	3.20
1327	Maruti	Hyderabad	2015	50205	Petrol	Manual	1	16.10	NaN	NaN	NaN	NaN	5.80
1385	Land	Pune	2014	11000	Petrol	Manual	2	18.15	NaN	NaN	NaN	NaN	1.90
1400	Land	Coimbatore	2008	69076	Petrol	Manual	1	18.15	NaN	NaN	NaN	NaN	40.86
2074	Maruti	Pune	2011	24255	Petrol	Manual	1	16.10	NaN	NaN	NaN	NaN	3.15
2096	Hyundai	Coimbatore	2004	52145	Petrol	Manual	1	18.15	NaN	NaN	NaN	NaN	1.93
2264	Toyota	Pune	2012	24500	Petrol	Manual	2	18.30	NaN	NaN	NaN	NaN	2.95
2325	Maruti	Pune	2015	67000	Petrol	Manual	1	16.10	NaN	NaN	NaN	NaN	1.75
2335	Maruti	Mumbai	2007	55000	Petrol	Manual	2	16.10	NaN	NaN	NaN	NaN	1.75
2530	BMW	Kochi	2014	64158	Diesel	Automatic	1	18.48	NaN	NaN	NaN	NaN	17.89
2542	Hyundai	Bangalore	2011	65000	Petrol	Manual	2	18.15	NaN	NaN	NaN	NaN	3.15
2623	BMW	Pune	2012	95000	Diesel	Automatic	2	16.48	NaN	NaN	NaN	NaN	18.90
2658	Maruti	Kolkata	2014	32986	Petrol	Manual	1	16.10	NaN	NaN	NaN	NaN	4.34
2737	Maruti	Japur	2001	20000	Petrol	Manual	1	12.90	NaN	NaN	NaN	NaN	0.70
2780	Hyundai	Pune	2009	100000	Petrol	Manual	1	18.15	NaN	NaN	NaN	NaN	1.60
2842	Hyundai	Bangalore	2012	43000	Petrol	Manual	1	18.15	NaN	NaN	NaN	NaN	3.25
3272	BMW	Mumbai	2008	81000	Diesel	Automatic	2	18.48	NaN	NaN	NaN	NaN	10.50
3404	Maruti	Japur	2006	125000	Petrol	Manual	4	16.10	NaN	NaN	NaN	NaN	2.35
3520	BMW	Delhi	2012	90000	Diesel	Automatic	1	18.48	NaN	NaN	NaN	NaN	14.50
3522	Hyundai	Kochi	2012	66400	Petrol	Manual	1	18.15	NaN	NaN	NaN	NaN	2.66
3830	Honda	Kolkata	2013	27000	Petrol	Automatic	1	14.00	NaN	NaN	NaN	NaN	11.99
4185	Land	Pune	2011	4927	Diesel	Manual	1	20.80	NaN	NaN	NaN	NaN	2.60
4182	Land	Mumbai	2003	75000	Diesel	Automatic	2	18.15	NaN	NaN	NaN	NaN	16.11
4229	Hyundai	Bangalore	2005	79000	Petrol	Manual	2	17.00	NaN	NaN	NaN	NaN	1.65
4577	BMW	Delhi	2012	72000	Diesel	Automatic	3	18.48	NaN	NaN	NaN	NaN	13.85
4604	Honda	Pune	2011	88000	Petrol	Manual	1	16.70	NaN	NaN	NaN	NaN	3.15
4697	Fiat	Kochi	2017	17941	Petrol	Manual	1	15.70	NaN	NaN	NaN	NaN	3.93
4712	Hyundai	Pune	2003	80000	Petrol	Manual	2	17.00	NaN	NaN	NaN	NaN	0.90
4952	Fiat	Kolkata	2010	47000	Petrol	Manual	1	14.60	NaN	NaN	NaN	NaN	1.49
5015	Maruti	Delhi	2006	63000	Petrol	Manual	1	16.10	NaN	NaN	NaN	NaN	1.60
5094	Hyundai	Delhi	2012	52000	Petrol	Manual	1	16.10	NaN	NaN	NaN	NaN	3.05
5270	Honda	Bangalore	2002	53000	Petrol	Manual	2	18.15	NaN	NaN	NaN	NaN	1.85

```
In [10]: train['Power'] = train['Power'].str.extract('([0-9]+\.?[0-9]+)')
train['Power'] = train['Power'].astype(float)
train['Engine'] = train['Engine'].str.extract('([0-9]+)')
train['Engine'] = train['Engine'].astype(float)

#Power의 Engine 모두 숫자였기때문에 남기고 float형으로 바꾼다.

In [11]: train_num = train.drop(['Name', 'Location', 'Fuel_Type', 'Transmission', 'New_Price'], axis = 1)
#숫자만 있는 데이터로 모델링 할것인 경우, 행의 인덱스가 필요치않아,
#데이터를 다루기에는, 특정한 위치의 결측값을 문자로 할 것이기때문에 열은 모두 지운다.
#drop()함수를 사용한다. axis = 1은 열을 지운다는 뜻이다.
```

Year, Kilometers_Driven

```
In [12]: train['Kilometers_Driven'] = train['Kilometers_Driven'].astype(float)

In [13]: from sklearn.impute import SimpleImputer
transformer = SimpleImputer()
transformer.fit(train_num)
x = transformer.transform(train_num)
#SimpleImputer의 transformer를 이용하여 결측치를 채운다.
train_num = pd.DataFrame(x, columns = train_num.columns, index = train_num.index)

In [14]: import matplotlib.pyplot as plt
plt.figure(figsize=(20,10))
for i in range(0,6):
    plt.subplot(2,3,i+1)
    plt.scatter(train_num.iloc[:,i], train_num['Price'])
    plt.title(train_num.columns[i])
#train 데이터의 숫자형데이터들을 matplotlib으로 시각화
#train 데이터의 숫자형 데이터를 matplotlib로 시각화
```



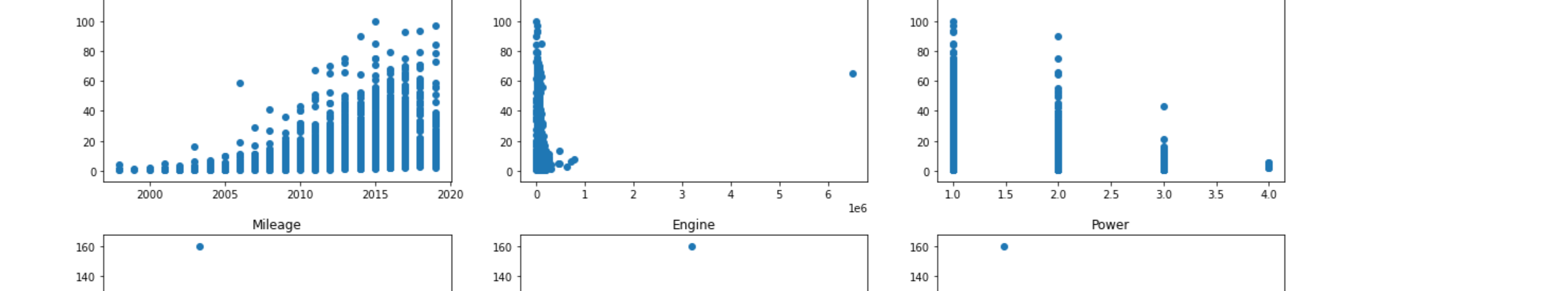
```
In [15]: train_num[train_num['Kilometers_Driven'] > 6500000]
#2328, 2328의 'Kilometers_Driven' 값이 잘못 입력된 것 같다.
#데이터가 매우 이상치로, 데이터를 삭제하고 잘 생각해 보자.
train_num.drop(2328, inplace = True)

In [16]: train_num[train_num['Kilometers_Driven'] > 6500000]

Out[16]:
```

Year	Kilometers_Driven	Owner_Type	Mileage	Engine	Power	Seats	Price
2000	2000	1	2000	2000	2000	2000	2000

```
In [17]: import matplotlib.pyplot as plt
plt.figure(figsize=(20,10))
for i in range(0,6):
    plt.subplot(2,3,i+1)
    plt.scatter(train_num.iloc[:,i], train_num['Price'])
    plt.title(train_num.columns[i])
#train 데이터의 숫자형데이터들을 matplotlib으로 시각화
#train 데이터의 숫자형 데이터를 matplotlib로 시각화
```



```
In [18]: train_cat = train.drop(['Year', 'Kilometers_Driven', 'Owner_Type', 'Mileage', 'Engine', 'Power', 'Seats', 'New_Price', 'Price'], axis = 1)
train_num = train.drop(['Year', 'Kilometers_Driven', 'Owner_Type', 'Mileage', 'Engine', 'Power', 'Seats', 'New_Price', 'Price'], axis = 1)
#연주력 데이터들을 처리하기 위해, train_cat을 생성한다.
#train_cat은 train_num에서와 같은 구조이지만, 행의 인덱스가 2328로 채워져 있다.
train_cat = train_cat.drop(2328, inplace = True)

In [19]: train_cat

Out[19]:
```

	Name	Location	Fuel_Type	Transmission
0	Maruti	Mumbai	CNG	Manual
1	Hyundai	Pune	Diesel	Manual
2	Honda	Chennai	Petrol	Manual
3	Maruti	Chennai	Diesel	Manual
4	Audi	Coimbatore	Diesel	Automatic
...
6014	Maruti	Delhi	Diesel	Manual
6015	Hyundai	Japur	Diesel	Manual
6016	Mahindra	Japur	Diesel	Manual
6017	Maruti	Kolkata	Petrol	Manual
6018	Chevrolet	Hyderabad	Diesel	Manual

6018 rows x 4 columns

```
In [20]: name = pd.DataFrame(train_cat.iloc[:,0])
loc = pd.DataFrame(train_cat.iloc[:,1])
fuel = pd.DataFrame(train_cat.iloc[:,2])
trans = pd.DataFrame(train_cat.iloc[:,3])
#iloc[]으로 접근해서, 인덱스 값은 주지 않아도 된다.
#name, loc, fuel, trans를 train_num에 추가한다.
train_num = pd.concat([name, loc, fuel, trans], axis = 1)
#name, loc, fuel, trans를 train_num에 추가한다.
train_cat = train_cat.drop(['Name', 'Location', 'Fuel_Type', 'Transmission'], axis = 1)

In [21]: train_prepared = pd.concat([train_scaled, train_cat], axis = 1)
train_labels = train_scaled['Price']

#train_scaled, train_cat, train_prepared을 생성한다.
#train_prepared은 train_labels과 같은 구조이지만, 행의 인덱스가 2328로 채워져 있다.
#train_prepared에 더할 데이터를 지정한다.

In [22]: train_prepared[train_prepared['Kolkata'] != 0]
#2328이 존재하지 않는 데이터는 train_labels과 같은 구조이지만, 행의 인덱스가 2328로 채워져 있다.
```

```
In [23]: train_prepared

Out[23]:
```

	Year	Kilometers_Driven	Owner_Type	Mileage	Engine	Power	Seats	Price	Name_Audi	Name_BMW	...	Location_Japur	Location_Kochi	Location_Kolkata	Location_Mumbai	Location_Pune	Fuel_Type
0	1.02																