

Cleaning Dataset

In [5]: `car.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 892 entries, 0 to 891
Data columns (total 6 columns):
#   Column          Non-Null Count  Dtype
---  -
0   name            892 non-null   object
1   company         892 non-null   object
2   year            892 non-null   object
3   Price           892 non-null   object
4   kms_driven      840 non-null   object
5   fuel_type       837 non-null   object
dtypes: object(6)
memory usage: 41.9+ KB
```

Quality

- ⌚ names are pretty inconsistent
- ⌚ names have company names attached to it
- ⌚ some names are spam like 'Maruti Ertiga showroom condition with' and 'Well maintained Tata Sumo'
- ⌚ company: many of the names are not of any company like 'Used', 'URJENT', and so on.
- ⌚ year has many non-year values
- ⌚ year is in object. Change to integer
- ⌚ Price has Ask for Price
- ⌚ Price has commas in its prices and is in object
- ⌚ kms_driven has object values with kms at last.
- ⌚ It has nan values and two rows have 'Petrol' in them
- ⌚ fuel_type has nan values

Cleaning Data

year has many non-year values

```
In [7]: car=car[car['year'].str.isnumeric()]
```

year is in object. Change to integer

```
In [8]: car['year']=car['year'].astype(int)
```

Price has Ask for Price

```
In [9]: car=car[car['Price']!='Ask For Price']
```

Price has commas in its prices and is in object

```
In [10]: car['Price']=car['Price'].str.replace(',','').astype(int)
```

kms_driven has object values with kms at last.

```
In [11]: car['kms_driven']=car['kms_driven'].str.split().str.get(0).str.replace(' ','')
```

It has nan values and two rows have 'Petrol' in them

```
In [12]: car=car[car['kms_driven'].str.isnumeric()]
```

```
In [13]: car['kms_driven']=car['kms_driven'].astype(int)
```

fuel_type has nan values

```
In [14]: car=car[~car['fuel_type'].isna()]
```

```
In [15]: car.shape
```

```
Out[15]: (816, 6)
```

name and **company** had spammed data...but with the previous cleaning, those rows got removed.

Company does not need any cleaning now. Changing car names. Keeping only the first three words

```
In [16]: car['name']=car['name'].str.split().str.slice(start=0,stop=3).str.join(' ')
```

Resetting the index of the final cleaned data

```
In [17]: car=car.reset_index(drop=True)
```

Cleaned Data

```
In [18]: car
```

```
Out[18]:
```

	name	company	year	Price	kms_driven	fuel_type
0	Hyundai Santro Xing	Hyundai	2007	80000	45000	Petrol
1	Mahindra Jeep CL550	Mahindra	2006	425000	40	Diesel
2	Hyundai Grand i10	Hyundai	2014	325000	28000	Petrol
3	Ford EcoSport Titanium	Ford	2014	575000	36000	Diesel
4	Ford Figo	Ford	2012	175000	41000	Diesel
...
811	Maruti Suzuki Ritz	Maruti	2011	270000	50000	Petrol
812	Tata Indica V2	Tata	2009	110000	30000	Diesel
813	Toyota Corolla Altis	Toyota	2009	300000	132000	Petrol
814	Tata Zest XM	Tata	2018	260000	27000	Diesel
815	Mahindra Quanto C8	Mahindra	2013	390000	40000	Diesel

816 rows × 6 columns

In [20]: `car.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 816 entries, 0 to 815
Data columns (total 6 columns):
#   Column      Non-Null Count  Dtype
---  -
0   name        816 non-null   object
1   company     816 non-null   object
2   year        816 non-null   int32
3   Price       816 non-null   int32
4   kms_driven  816 non-null   int32
5   fuel_type   816 non-null   object
dtypes: int32(3), object(3)
memory usage: 28.8+ KB
```
