

```
In [1]: import pandas as pd
import numpy as np
matplotlib inline
import matplotlib.pyplot as plt
```

```
In [2]: df = pd.read_csv('auto.csv')
df.head(5)
```

	3	7	alfa-romero	gas	std	two	convertible	rwd	front	88.6	...	130	mpfi	3.47	2.68	9.0	111	5000	21	27	13495
0	3	7	alfa-romero	gas	std	two	convertible	rwd	front	88.6	...	130	mpfi	3.47	2.68	9.0	111	5000	21	27	16500
1	1	7	alfa-romero	gas	std	two	hatchback	rwd	front	94.5	...	152	mpfi	2.68	3.47	9.0	154	5000	19	26	16500
2	2	164	audi	gas	std	four	sedan	fwd	front	99.8	...	109	mpfi	3.19	3.40	10.0	102	5500	24	30	13950
3	2	164	audi	gas	std	four	sedan	4wd	front	99.4	...	136	mpfi	3.19	3.40	10.0	115	5500	18	22	17450
4	2	7	audi	gas	std	two	sedan	fwd	front	99.8	...	136	mpfi	3.19	3.40	8.5	110	5500	19	25	12500
...
199	-1	95	volvo	gas	std	four	sedan	rwd	front	109.1	...	141	mpfi	3.78	3.15	9.5	114	5400	23	28	16845
200	-1	95	volvo	gas	turbo	four	sedan	rwd	front	109.1	...	141	mpfi	3.78	3.15	8.7	160	5300	19	25	19045
201	-1	95	volvo	gas	std	four	sedan	rwd	front	109.1	...	173	mpfi	3.58	2.87	8.8	134	5500	18	23	21485
202	-1	95	volvo	diesel	turbo	four	sedan	rwd	front	109.1	...	145	kdi	3.01	3.40	23.0	106	4800	26	27	22470
203	-1	95	volvo	gas	turbo	four	sedan	rwd	front	109.1	...	141	mpfi	3.78	3.15	9.5	114	5400	19	25	22625

204 rows × 26 columns

```
In [3]: # Let's give header to the datasets
headers = ["symboling", "normalized-losses", "make", "fuel-type", "aspiration", "num-of-doors", "body-style",
           "drive-wheels", "engine-location", "wheel-base", "length", "width", "height", "curb-weight", "engine-type",
           "num-of-cylinders", "engine-size", "fuel-system", "bore", "stroke", "compression-ratio", "horsepower",
           "peak-rpm", "city-mpg", "highway-mpg", "price"]
```

```
In [4]: df.columns=headers
```

```
In [5]: df.head(5)
```

```
In [6]: # Convert '' to NaN
df.replace('','',np.nan, inplace=True)
df.head(5)
```

	symboling	normalized-losses	make	fuel-type	aspiration	num-of-doors	body-style	drive-wheels	engine-location	wheel-base	engine-size	fuel-system	bore	stroke	compression-ratio	horsepower	peak-rpm	city-mpg	highway-mpg	price	
0	3	NaN	alfa-romero	gas	std	two	convertible	rwd	front	88.6	...	130	mpfi	3.47	2.68	9.0	111	5000	21	27	13495
1	1	NaN	alfa-romero	gas	std	two	hatchback	rwd	front	94.5	...	152	mpfi	2.68	3.47	9.0	154	5000	19	26	16500
2	2	164	audi	gas	std	four	sedan	fwd	front	99.8	...	109	mpfi	3.19	3.40	10.0	102	5500	24	30	13950
3	2	164	audi	gas	std	four	sedan	4wd	front	99.4	...	136	mpfi	3.19	3.40	10.0	115	5500	18	22	17450
4	2	NaN	audi	gas	std	two	sedan	fwd	front	99.8	...	136	mpfi	3.19	3.40	8.5	110	5500	19	25	12500

5 rows × 26 columns

```
In [7]: # Checking the null value in each column
df.isnull().sum()
```

symboling	0
normalized-losses	40
make	0
fuel-type	0
aspiration	0
num-of-doors	2
body-style	0
drive-wheels	0
engine-location	0
wheel-base	0
length	0

```
Out[7]:
```

Identifying and handling missing data

```
In [6]: # Convert 'NaN' to NaN
df.replace('?', np.nan, inplace=True)
df.head(5)
```

```
highway-mpg      0
price            4
dtype: int64
```

Handling this missing value based on their characteristics.

```
In [8]: #Replacing "NaN" with mean value in "normalized-losses" column
avg_norm_loss = df["normalized-losses"].astype("float").mean(axis=0)
df["normalized-losses"].replace(np.nan, avg_norm_loss, inplace=True)

In [9]: #replace the missing 'num-of-doors' values by the most frequent
df["num-of-doors"].value_counts()

Out[9]:
four    114
two      88
Name: num-of-doors, dtype: int64

In [10]: df["num-of-doors"].replace(np.nan, "four", inplace=True)

In [11]: #Replacing "NaN" with mean value in "bore" column
avg_bore_df = df["bore"].astype("float").mean(axis=0)
df["bore"].replace(np.nan, avg_bore, inplace=True)
```

```
In [7]: # Checking the null value in each column
df.isnull().sum()
```

```
Out[7]:
```

symboling	0
normalized-losses	40
make	0
fuel-type	0
aspiration	0
num-of-doors	2
body-style	0
drive-wheels	0
engine-location	0
wheel-base	0
length	0
width	0
height	0
curb-weight	0
engine-type	0
num-of-cylinders	0
engine-size	0
fuel-system	0
bore	4
stroke	4
compression-ratio	0
horsepower	2
peak-rpm	2
city-mpg	0
highway-mpg	0
price	4
dtype:	int64

Handling this missing value based on their characteristics.

```
In [8]: #Replacing "NaN" with mean value in "normalized-losses" column
avg_norm_loss = df["normalized-losses"].astype("float").mean(axis=0)
df["normalized-losses"].replace(np.nan, avg_norm_loss, inplace=True)
```

```
In [9]: #Replace the missing 'num-of-doors' values by the most frequent
df["num-of-doors"].value_counts()
```

```
Out[9]:
```

four	11
two	88

Name: num-of-doors, dtype: int64

```
In [10]: df["num-of-doors"].replace(np.nan, "four", inplace=True)
```

```
In [11]: #Replacing "NaN" with mean value in "bore" column
avg_bore=df['bore'].astype('float').mean(axis=0)
df["bore"].replace(np.nan, avg_bore, inplace=True)
```

```
In [12]: #Replacing "NaN" with mean value in "stroke" column
avg_stroke=df['stroke'].astype('float').mean(axis=0)
df["stroke"].replace(np.nan, avg_stroke, inplace=True)
```

```
In [13]: #Replacing "NaN" with the mean value in the "horsepower" column
avg_horsepower = df['horsepower'].astype('float').mean(axis=0)
df["horsepower"].replace(np.nan, avg_horsepower, inplace=True)
```

```
In [14]: #Replacing "NaN" with the mean value in the "peak-rpm" column
avg_peakrpm=df['peak-rpm'].astype('float').mean(axis=0)
df["peak-rpm"].replace(np.nan, avg_peakrpm, inplace=True)
```

```
In [15]: #Finally, let's drop all rows that do not have price data:
df.dropna(subset=["price"], axis=0, inplace=True)
```

```
# reset index, because we dropped two rows
df.reset_index(drop=True, inplace=True)
```

```
In [16]: df.head()
```

```
aspiration      object
num-of-doors   object
body-style      object
drive-wheels    object
engine-location object
wheel-base     float64
length         float64
width          float64
height         float64
curb-weight     int64
engine-type     object
num-of-cylinders object
engine-size     int64
fuel-system     object
bore           float64
stroke         float64
compression-ratio float64
horsepower      int32
peak-rpm       float64
city-mpg        int64
highway-mpg     int64
price          float64
dtype: object
```

Let's standardized and normalized the data

```
In [20]: # transform mpg to L/100km by mathematical operation (235 divided by mpg)
df["highway-mpg"] = 235/df["highway-mpg"]

# rename column name from "highway-mpg" to "highway-L/100km"
df.rename(columns={"highway-mpg": "highway-L/100km"}, inplace=True)

In [21]: # transform mpg to L/100km by mathematical operation (235 divided by mpg)
df["city-mpg"] = 235/df["city-mpg"]
```

```
In [17]: #Let's list the data types for each column
df.dtypes
```

```
Out[17]:
```

symboling	int64
normalized-losses	object
make	object
fuel-type	object
aspiration	object
num-of-doors	object
body-style	object
drive-wheels	object
engine-location	object
wheel-base	float64
length	float64
width	float64
height	float64
curb-weight	int64
engine-type	object
num-of-cylinders	int64
engine-size	int64
fuel-system	object
bore	object
stroke	object
compression-ratio	float64
horsepower	object
peak-rpm	object
city-mpg	int64
highway-mpg	int64
price	object
dtype:	object

```
In [18]: #Convert data types to proper format
df["bore"] = df["bore"].astype("float")
df["stroke"] = df["stroke"].astype("float")
df["normalized-losses"] = df["normalized-losses"].astype("int")
df["price"] = df["price"].astype("float")
df["peak-rpm"] = df["peak-rpm"].astype("int")
df["horsepower"] = df["horsepower"].astype("int")
```

```
In [19]: df.dtypes
```

```
Out[19]:
```

symboling	int64
normalized-losses	int32
make	object
fuel-type	object
aspiration	object
num-of-doors	object
body-style	object
drive-wheels	object
engine-location	object
wheel-base	float64
length	float64
width	float64
height	float64
curb-weight	int64
engine-type	object
num-of-cylinders	object
engine-size	int64
fuel-system	object
bore	float64
stroke	float64
compression-ratio	float64
horsepower	int32
peak-rpm	float64
city-mpg	int64
highway-mpg	int64
price	float64
dtype:	object

Let's standardized and normalized the data

```
In [20]: # transform mpg to 1/100km by mathematical operation (235 divided by mpg)
df["highway-mpg"] = 235/df["highway-mpg"]
```

```
# rename column name from "highway-mpg" to "highway-L/100km"
df.rename(columns={"highway-mpg": "highway-L/100km", inplace=True)
```

```
In [21]: # transform mpg to 1/100km by mathematical operation (235 divided by mpg)
df["city-mpg"] = 235/df["city-mpg"]
```

```
# rename column name from "city-mpg" to "city-L/100km"
df.rename(columns={"city-mpg": "city-L/100km", inplace=True)
```

```
In [22]: df.head()
```

```
Out[22]:
```

	symboling	normalized-losses	make	fuel-type	aspiration	num-of-doors	body-style	drive-wheels	engine-location	wheel-base	engine-size	fuel-system	bore	stroke	compression-ratio	horsepower	peak-rpm	city-L/100km	highway-L/100km	price
0	3	122	alfa-romero	gas	std	two	convertible	rwd	front	88.6	...	130	mpfi	3.47	2.68	9.0	111	19.3	27.7	13495
1	1	122	alfa-romero	gas	std	two	hatchback	rwd	front	94.5	...	152	mpfi	2.68	3.47	9.0	154	15.5	26.3	16500
2	2	164	audi	gas	std	four	sedan	fwd	front	99.8	...	109	mpfi	3.19	3.40	10.0	102	23.3	30.0	13950
3	2	164	audi	gas	std	four	sedan	4wd	front	99.4	...	136	mpfi	3.19	3.40	10.0	115	17.0	22.7	17450
4	2	122	audi	gas	std	two	sedan	fwd	front	99.8	...	136	mpfi	3.19	3.40	8.5	110	19.4	25.0	12500

5 rows × 26 columns

```
In [23]: #I would like to normalize some variables so their value ranges from 0 to 1
# replace (original value) by (original value)/(maximum value)
df['length'] = df['length']/df['length'].max()
df['width'] = df['width']/df['width'].max()
df['height'] = df['height']/df['height'].max()
```

```
In [24]: #Let's plot the histogram of horsepower to see what the distribution of horsepower looks like.
plt.hist(df['horsepower'])
plt.xlabel("horsepower")
plt.ylabel("count")
plt.title("horsepower")
```

```
Out[24]:
```

Text(0.5, 1.0, 'horsepower')



```
In [25]: bins = np.linspace(min(df["horsepower"]), max(df["horsepower"]), 4)
group_names = ['Low', 'Medium', 'High']
df["horsepower-binned"] = pd.cut(df["horsepower"], bins, labels=group_names, include_lowest=True)
df["horsepower-binned"].head(10)
```

```
Out[25]:
```

	horsepower	horsepower-binned
0	111	Low
1	154	Medium
2	102	Low
3	115	Low
4	110	Low
5	110	Low
6	110	Low
7	140	Medium
8	101	Low
9	101	Low

```
In [26]: df["horsepower-binned"].value_counts()
```

```
Out[26]:
```

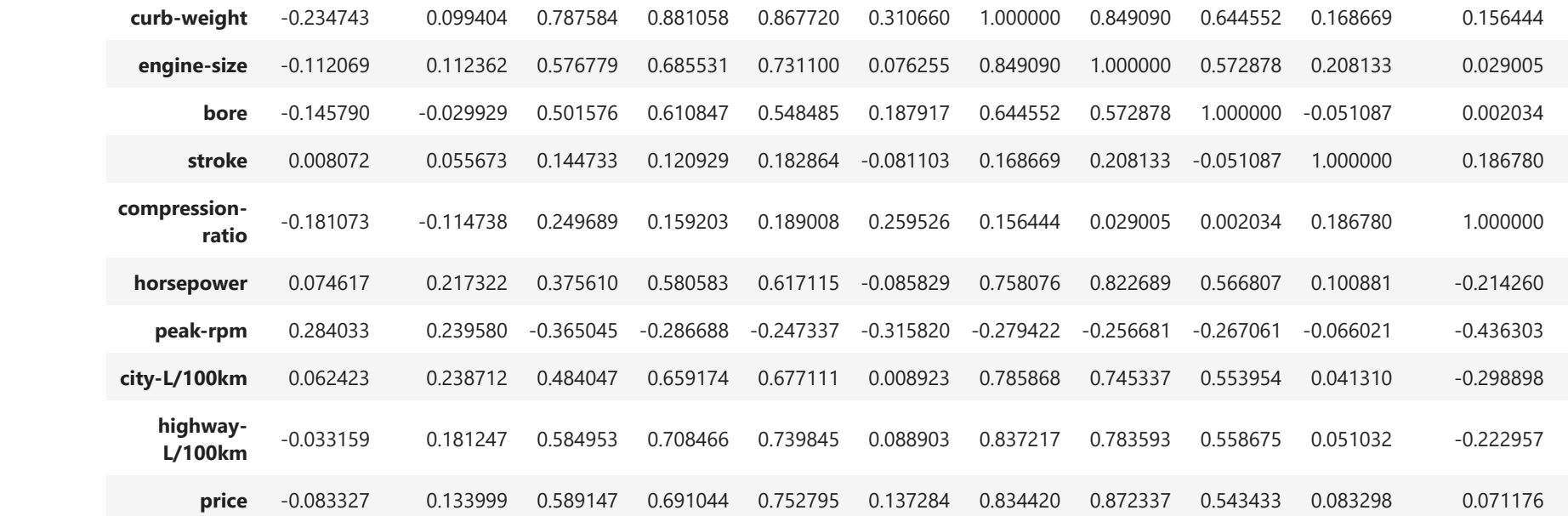
Low	152
Medium	43
High	5

Name: horsepower-binned, dtype: int64

```
In [27]: plt.bar(group_names, df["horsepower-binned"].value_counts())
plt.xlabel("horsepower")
plt.ylabel("count")
plt.title("horsepower bins")
```

```
Out[27]:
```

Text(0.5, 1.0, 'horsepower bins')



```
In [28]: #We use indicator variables so we can use categorical variables for regression analysis later.
dummy_variable_1 = pd.get_dummies(df["fuel-type"])
dummy_variable_1.head()
```

```
Out[28]:
```

	diesel	gas
0	0	1
1	0	1
2	0	1
3	0	1
4	0	1

```
In [29]: dummy_variable_1.rename(columns={'gas':'fuel-type-gas', 'diesel':'fuel-type-diesel', inplace=True)
dummy_variable_1.head()
```

```
Out[29]:
```

	fuel-type-diesel	fuel-type-gas
0	0	1
1	0	1
2	0	1
3	0	1
4	0	1

```
In [30]: # merge data frame "df" and "dummy_variable_1"
df = pd.concat([df, dummy_variable_1], axis=1)
```

```
# drop original column "fuel-type" from "df"
df.drop("fuel-type", axis = 1, inplace=True)
```

```
Out[31]: df.head()
```

```
Out[31]:
```

	symboling	normalized-losses	make	fuel-type	aspiration	num-of-doors	body-style	drive-wheels	engine-location	wheel-base	length	...	stroke	compression-ratio	horsepower
0	3	122	alfa-romero	gas	std	two	convertible	rwd	front	88.6	0.811148	...	2.68	9.0	111
1	1	122	alfa-romero	gas	std	two	hatchback	rwd	front	94.5	0.822681	...	3.47	9.0	154
2	2	164	audi	gas	std	four	sedan	fwd	front	99.8	0.848630	...	3.40	10.0	102
3	2	164	audi	gas	std	four	sedan	4wd	front	99.4	0.848630	...	3.40	10.0	115
4	2	122	audi	gas	std	two	sedan	fwd	front	99.8	0.851994	...	3.40	8.5	110

5 rows × 28 columns

```
In [32]: dummy_variable_2 = pd.get_dummies(df["aspiration"])
```

```
# change column names for clarity
dummy_variable_2.rename(columns={'std':'aspiration-std', 'turbo': 'aspiration-turbo', inplace=True)
# show first 5 instances of data frame "dummy_variable_1"
```

```
Out[32]: dummy_variable_2.head()
```

```
Out[32]:
```

	aspiration-std	aspiration-turbo
0	1	0
1	1	0
2	1	0

