

Data Analysis for Automobile car dealership with Python

```
1. Exploring the data for Automobile car dealership

In [2]: # Import pandas library
import pandas as pd

In [3]: #The online file by the URL provides above ,and assign it to variable "df"
other_path = "https://s3-ap-southeast-1.amazonaws.com/courses-data/CognitiveClass/DAB181EN/autos.csv"
df = pd.read_csv(other_path, header=None)

In [4]: # Showing the summary of the first 5 row using dataframe.head()
print("The first 5 rows of the dataframe")
df.head(5)

The first 5 rows of the dataframe
Out[4]:
   0    1    2    3    4    5    6    7    8    9  ...  16  17  18  19  20  21  22  23  24  25
0    3    ?  alfa-romeo  gas  std  two  convertible  fwd  front  88.6  ...  130  mpg  3.47  2.68  9.0  111  5000  21  27  13495
1    3    ?  alfa-romeo  gas  std  two  convertible  fwd  front  88.6  ...  130  mpg  3.47  2.68  9.0  111  5000  21  27  16500
2    1    ?  alfa-romeo  gas  std  two  hatchback  fwd  front  84.5  ...  152  mpg  2.68  3.47  9.0  154  5000  19  26  16500
3    2    164  audi  gas  std  four  sedan  fwd  front  99.8  ...  109  mpg  3.19  3.40  10.0  102  5500  24  30  13950
4    2    164  audi  gas  std  four  sedan  fwd  front  99.4  ...  136  mpg  3.19  3.40  8.0  115  5500  18  22  17450
5 rows x 26 columns

In [5]: #Showing the bottom 10 rows
df.tail(10)

Out[5]:
   196  1  74  volvo  gas  std  four  wagon  fwd  front  104.3  ...  141  mpg  3.78  3.15  9.5  114  5400  23  28  13415
197  -2  103  volvo  gas  std  four  sedan  fwd  front  104.3  ...  141  mpg  3.78  3.15  9.5  114  5400  24  28  15985
198  -1  74  volvo  gas  std  four  wagon  fwd  front  104.3  ...  141  mpg  3.78  3.15  9.5  114  5400  24  28  16515
199  -2  103  volvo  gas  turbo  four  sedan  fwd  front  104.3  ...  130  mpg  3.62  3.15  7.5  162  5100  17  22  18420
200  -1  74  volvo  gas  turbo  four  wagon  fwd  front  104.3  ...  130  mpg  3.62  3.15  7.5  162  5100  17  22  18950
201  -1  95  volvo  gas  std  four  sedan  fwd  front  109.1  ...  141  mpg  3.78  3.15  9.5  114  5400  23  28  16845
202  -1  95  volvo  gas  turbo  four  sedan  fwd  front  109.1  ...  141  mpg  3.78  3.15  8.7  160  5300  19  25  19045
203  -1  95  volvo  gas  std  four  sedan  fwd  front  109.1  ...  173  mpg  3.19  3.40  8.0  154  5000  19  26  16500
204  -3  95  volvo  diesel  turbo  four  sedan  fwd  front  109.1  ...  145  hp  3.01  3.40  23.0  106  4800  26  27  22470
205  -1  95  volvo  gas  turbo  four  sedan  fwd  front  109.1  ...  141  mpg  3.78  3.15  9.5  114  5400  19  25  22625
10 rows x 26 columns

In [7]: # creating headers for the rows.
headers = ["symboling","normalized-loses","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-rate","horsepower","peak-rpm","city-mpg","highway-mpg","price"]
print(headers)

headers
Out[8]:
symboling normalized-loses 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-rate horsepower peak-rpm city-mpg highway-mpg price

In [8]: # I will replace headers and recheck the data frame
df.columns=headers

Out[8]:
symboling normalized-loses make fuel-type aspiration num-of-doors body-style drive-wheels engine-location wheel-base ... engine-size fuel-system bore stroke compression-rate horsepower peak-rpm city-mpg highway-mpg price
0 3 ? alfa-romeo gas std two convertible fwd front 88.6 ... 130 mpg 3.47 2.68 9.0 111 5000 21 27 13495
1 3 ? alfa-romeo gas std two convertible fwd front 88.6 ... 130 mpg 3.47 2.68 9.0 111 5000 21 27 16500
2 1 ? alfa-romeo gas std two hatchback fwd front 84.5 ... 152 mpg 2.68 3.47 9.0 154 5000 19 26 16500
3 2 164 audi gas std four sedan fwd front 99.8 ... 109 mpg 3.19 3.40 10.0 102 5500 24 30 13950
4 2 164 audi gas std four sedan fwd front 99.4 ... 136 mpg 3.19 3.40 8.0 115 5500 18 22 17450
5 2 ? audi gas std four sedan fwd front 99.8 ... 136 mpg 3.19 3.40 8.5 110 5500 19 25 16250
6 1 150 audi gas std four sedan fwd front 105.8 ... 126 mpg 3.19 3.40 8.5 110 5500 19 26 16500
7 1 ? audi gas std four wagon fwd front 105.8 ... 136 mpg 3.19 3.40 8.5 110 5500 19 25 17710
8 1 150 audi gas turbo four sedan fwd front 105.8 ... 131 mpg 3.13 3.40 8.3 140 5500 17 20 23875
9 0 ? audi gas turbo two hatchback fwd front 99.5 ... 131 mpg 3.13 3.40 7.0 160 5500 16 22 7
10 rows x 26 columns

In [9]: # Dropping missing values along the column "price" so that i can focus on the cars that were bought.
df.dropna(subset=["price"], axis=0)

Out[9]:
symboling normalized-loses make fuel-type aspiration num-of-doors body-style drive-wheels engine-location wheel-base ... engine-size fuel-system bore stroke compression-rate horsepower peak-rpm city-mpg highway-mpg price
0 3 ? alfa-romeo gas std two convertible fwd front 88.6 ... 130 mpg 3.47 2.68 9.0 111 5000 21 27 13495
1 3 ? alfa-romeo gas std two convertible fwd front 88.6 ... 130 mpg 3.47 2.68 9.0 111 5000 21 27 16500
2 1 ? alfa-romeo gas std two hatchback fwd front 84.5 ... 152 mpg 2.68 3.47 9.0 154 5000 19 26 16500
3 2 164 audi gas std four sedan fwd front 99.8 ... 109 mpg 3.19 3.40 10.0 102 5500 24 30 13950
4 2 164 audi gas std four sedan fwd front 99.4 ... 136 mpg 3.19 3.40 8.0 115 5500 18 22 17450
... ..
200 -1 95 volvo gas std four sedan fwd front 109.1 ... 141 mpg 3.78 3.15 9.5 114 5400 23 28 16845
201 -1 95 volvo gas turbo four sedan fwd front 109.1 ... 141 mpg 3.78 3.15 8.7 160 5300 19 25 19045
202 -1 95 volvo gas std four sedan fwd front 109.1 ... 173 mpg 3.58 2.87 8.8 164 5500 18 23 21485
203 -1 95 volvo diesel turbo four sedan fwd front 109.1 ... 145 hp 3.01 3.40 23.0 106 4800 26 27 22470
204 -1 95 volvo gas turbo four sedan fwd front 109.1 ... 141 mpg 3.78 3.15 9.5 114 5400 19 25 22625
205 rows x 26 columns

In [14]: # Saving the Data set.
df.to_csv("automobile.csv", index=False)

2.Data Wrangling to convert the data from the initial format to a format that may be better for analysis.

In [15]: # libraries
import pandas, as pd
import matplotlib.pyplot as plt
import numpy as np

In [16]: # replace "?" to NaN
df.replace("?",np.nan , inplace =True)
df.head(5)

Out[16]:
symboling normalized-loses make fuel-type aspiration num-of-doors body-style drive-wheels engine-location wheel-base ... engine-size fuel-system bore stroke compression-rate horsepower peak-rpm city-mpg highway-mpg price
0 3 NaN alfa-romeo gas std two convertible fwd front 88.6 ... 130 mpg 3.47 2.68 9.0 111 5000 21 27 13495
1 3 ? alfa-romeo gas std two convertible fwd front 88.6 ... 130 mpg 3.47 2.68 9.0 111 5000 21 27 16500
2 1 ? alfa-romeo gas std two hatchback fwd front 84.5 ... 152 mpg 2.68 3.47 9.0 154 5000 19 26 16500
3 2 164 audi gas std four sedan fwd front 99.8 ... 109 mpg 3.19 3.40 10.0 102 5500 24 30 13950
4 2 164 audi gas std four sedan fwd front 99.4 ... 136 mpg 3.19 3.40 8.0 115 5500 18 22 17450
5 rows x 26 columns

In [17]: # Identifying missing values within the data.
missing_data = df.isnull()
missing_data.head(5)

Out[17]:
symboling normalized-loses make fuel-type aspiration num-of-doors body-style drive-wheels engine-location wheel-base ... engine-size fuel-system bore stroke compression-rate horsepower peak-rpm city-mpg highway-mpg price
0 False True False False False False False False False False False False False False False False
1 False True False False False False False False False False False False False False False False
2 False True False False False False False False False False False False False False False False
3 False False False False False False False False False False False False False False False False
4 False False False False False False False False False False False False False False False False
5 rows x 26 columns

In [18]: avg_norm_loss = df["normalized-loses"].astype("float").mean(axis=0)
print("Average of normalized-loses:",avg_norm_loss)

Average of normalized-loses: 122.8

In [19]: # COUNT MISSING VALUES IN EACH COLUMN
for column in missing_data.columns.values.tolist():
    print(column)
    print(missing_data[column].value_counts())
    print("")

symboling
False 285
Name: symboling, dtype: int64

normalized-loses
False 164
True 41
Name: normalized-loses, dtype: int64

make
False 285
Name: make, dtype: int64

fuel-type
False 285
Name: fuel-type, dtype: int64

aspiration
False 285
Name: aspiration, dtype: int64

num-of-doors
False 283
True 2
Name: num-of-doors, dtype: int64

body-style
False 285
Name: body-style, dtype: int64

drive-wheels
False 285
Name: drive-wheels, dtype: int64

engine-location
False 285
Name: engine-location, dtype: int64

wheel-base
False 285
Name: wheel_base, dtype: int64

length
False 285
Name: length, dtype: int64

width
False 285
Name: width, dtype: int64

height
False 285
Name: height, dtype: int64

curb-weight
False 285
Name: curb-weight, dtype: int64

engine-type
False 285
Name: engine-type, dtype: int64

num-of-cylinders
False 285
Name: num-of-cylinders, dtype: int64

engine-size
False 285
Name: engine-size, dtype: int64

fuel-system
False 285
Name: fuel-system, dtype: int64

bore
False 281
True 4
Name: bore, dtype: int64

stroke
False 281
True 4
Name: stroke, dtype: int64

compression-rate
False 285
Name: compression-rate, dtype: int64

horsepower
False 283
True 2
Name: horsepower, dtype: int64

peak-rpm
False 283
True 2
Name: peak-rpm, dtype: int64

city-mpg
False 285
Name: city-mpg, dtype: int64

highway-mpg
False 285
Name: highway-mpg, dtype: int64

price
False 281
True 4
Name: price, dtype: int64

3.Replacing missing values/null values on the data frame

In [20]: # Calculating average for normalized loses
avg_norm_loss = df["normalized-loses"].astype("float").mean(axis=0)
print("Average of normalized-loses: ",avg_norm_loss)

Average of normalized-loses: 122.8

In [21]: # Replacing missing value with normalized loses average on the data frame.
df["normalized-loses"].replace(np.nan,avg_norm_loss, inplace=True)

In [22]: #Calculating average for bore
avg_bore=df['bore'].astype('float').mean(axis=0)
print("Average of bore: ", avg_bore)

Average of bore: 3.3297512437810943

In [23]: #Replacng missing value bore average on the dataframe
df["bore"].replace(np.nan, avg_bore, inplace=True)

In [24]: #Calculating average for stroke
avg_stroke=df['stroke'].astype('float').mean(axis=0)
print("stroke: ", avg_stroke)

stroke: 3.295422885972139

In [25]: # Replacing missing value with stroke on the dataframe
df["stroke"].replace(np.nan,avg_stroke, inplace=True)

In [26]: # Calculating average for the horse power
avg_horsepower = df['horsepower'].astype('float').mean(axis=0)
print("Average of horsepower: ",avg_horsepower)

Average of horsepower: 184.25615783546799

In [27]: # Replacing missing value with horsepower on the dataframe.
df["horsepower"].replace(np.nan, avg_horsepower, inplace=True)

In [28]: # Calculating the average for peak-rpm
avg_peakrpm=df['peak-rpm'].astype('float').mean(axis=0)
print("peak-rpm: ", avg_peakrpm)

peak-rpm: 5125.369458128879

In [29]: #Replacng missing value peak-rpm on the dataframe
df["peak-rpm"].replace(np.nan, avg_peakrpm, inplace=True)

In [30]: # To see which values are present in a particular column
df["num-of-doors"].value_counts()

Out[30]:
four 114
two 89
Name: num-of-doors, dtype: int64

In [31]: # I can use the ".idxmax()" method to calculate the most comontype automatically:
df["num-of-doors"].value_counts().idxmax()

Out[31]:
'four'

In [32]: # Dropping the whole column with the missing data.
df.dropna(subset=["price"],axis=0, inplace=True)

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

In [33]: df.head()

Out[33]:
symboling normalized-loses make fuel-type aspiration num-of-doors body-style drive-wheels engine-location wheel-base ... engine-size fuel-system bore stroke compression-rate horsepower peak-rpm city-mpg highway-mpg price
0 3 122.0 alfa-romeo gas std two convertible fwd front 88.6 ... 130 mpg 3.47 2.68 9.0 111 5000.0 21 27 13495.0
1 3 122.0 alfa-romeo gas std two convertible fwd front 88.6 ... 130 mpg 3.47 2.68 9.0 111 5000.0 21 27 16500.0
2 1 122.0 alfa-romeo gas std two hatchback fwd front 84.5 ... 152 mpg 2.68 3.47 9.0 154 5000.0 19 26 16500.0
3 2 164 audi gas std four sedan fwd front 99.8 ... 109 mpg 3.19 3.40 10.0 102 5500.0 24 30 13950.0
4 2 164 audi gas std four sedan fwd front 99.4 ... 136 mpg 3.19 3.40 8.0 115 5500.0 18 22 17450.0
5 rows x 26 columns

1. Converting data types to proper formats

In [35]: # Converting data
df["bore","stroke"] = df[["bore","stroke"]].astype("float")
df["normalized-loses"] = df[["normalized-loses"]].astype("int")
df["price"] = df[["price"]].astype("float")
df[["peak-rpm"]] = df[["peak-rpm"]].astype("float")

In [37]: # list the columns after the conversion
df.dtypes

Out[37]:
symboling int64
normalized-loses 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 float64
fuel-system object
bore float64
stroke float64
compression-rate float64
horsepower object
peak-rpm float64
city-mpg int64
highway-mpg int64
price float64
dtype: object

1. Example

Transform miles per gallon to Kilometers.

In our dataset,the fuel consumption columns"city-mpg" and "highway-mpg" are represented by mpg(miles per gallon )unit.Assume we are developing an application in a country that accept the fuel consumption with L/100km standard.

We will need to apply data transformation to transform mpg into L/100km?

The formula for unit conversion is

L/100km =235/mpg

In [38]: df.head()

Out[38]:
symboling normalized-loses make fuel-type aspiration num-of-doors body-style drive-wheels engine-location wheel-base ... engine-size fuel-system bore stroke compression-rate horsepower peak-rpm city-mpg highway-mpg price
0 3 122 alfa-romeo gas std two convertible fwd front 88.6 ... 130 mpg 3.47 2.68 9.0 111 5000.0 21 27 13495.0
1 3 122 alfa-romeo gas std two convertible fwd front 88.6 ... 130 mpg 3.47 2.68 9.0 111 5000.0 21 27 16500.0
2 1 122 alfa-romeo gas std two hatchback fwd front 84.5 ... 152 mpg 2.68 3.47 9.0 154 5000.0 19 26 16500.0
3 2 164 audi gas std four sedan fwd front 99.8 ... 109 mpg 3.19 3.40 10.0 102 5500.0 24 30 13950.0
4 2 164 audi gas std four sedan fwd front 99.4 ... 136 mpg 3.19 3.40 8.0 115 5500.0 18 22 17450.0
5 rows x 26 columns

In [39]: # Convert mpg to L/100km by mathematical operation (235 divided by mpg)
df['city-L/100km'] = 235/df['city-mpg']

# check transformed data
df.head()

Out[39]:
symboling normalized-loses make fuel-type aspiration num-of-doors body-style drive-wheels engine-location wheel-base ... fuel-system bore stroke compression-rate horsepower peak-rpm city-mpg highway-mpg price city-L/100km
0 3 122 alfa-romeo gas std two convertible fwd front 88.6 ... mpg 3.47 2.68 9.0 111 5000.0 21 27 13495.0 11.90478
1 3 122 alfa-romeo gas std two convertible fwd front 88.6 ... mpg 3.47 2.68 9.0 111 5000.0 21 27 16500.0 11.90478
2 1 122 alfa-romeo gas std two hatchback fwd front 84.5 ... mpg 2.68 3.47 9.0 154 5000.0 19 26 16500.0 12.36421
3 2 164 audi gas std four sedan fwd front 99.8 ... mpg 3.19 3.40 10.0 102 5500.0 24 30 13950.0 9.79167
4 2 164 audi gas std four sedan fwd front 99.4 ... mpg 3.19 3.40 8.0 115 5500.0 18 22 17450.0 13.05556
5 rows x 27 columns

In [40]: # Transforming mpg to L/100km in the column of " highway-mpg "
df["highway-L/100km"] = 235/df["highway-mpg"]

# checkng transformed data
df.head()

Out[40]:
symboling normalized-loses make fuel-type aspiration num-of-doors body-style drive-wheels engine-location wheel-base ... bore stroke compression-rate horsepower peak-rpm city-mpg highway-mpg price city-L/100km highway-L/100km
0 3 122 alfa-romeo gas std two convertible fwd front 88.6 ... 3.47 2.68 9.0 111 5000.0 21 27 13495.0 11.90478 8.70374
1 3 122 alfa-romeo gas std two convertible fwd front 88.6 ... 3.47 2.68 9.0 111 5000.0 21 27 16500.0 11.90478 8.70374
2 1 122 alfa-romeo gas std two hatchback fwd front 84.5 ... 2.68 3.47 9.0 154 5000.0 19 26 16500.0 12.36421 9.03842
3 2 164 audi gas std four sedan fwd front 99.8 ... 3.19 3.40 10.0 102 5500.0 24 30 13950.0 9.79167 12.36421
4 2 164 audi gas std four sedan fwd front 99.4 ... 3.19 3.40 8.0 115 5500.0 18 22 17450.0 13.05556 10.68118
5 rows x 28 columns

1. Example: Objective I want to normalize the variables so their values ranges from 0 to 1. Approach: replace original value by (original value)/(maximum value)

In [41]: #replace(original value) by (original value)/(maximum value)
df["length"] = df["length"]/df["length"].max()
df["width"] = df["width"]/df["width"].max()

In [42]: # Normalizing the column "height"
df["height"] = df["height"]/df["height"].max()

1. Example in our dataset,"horsepower" is a real-valued variable ranging from 48 to 288.8,has 57 unique values.What if we only care about the price difference between cars with high horsepower,medium horsepower , and little (3)types?Can we rearrange them into three "bins" to simplify analysis?

I can only use the method "cut" to segment the "horsepower" column into 3 bins

In [44]: # Converting mpg to correct format
df["horsepower"]df["horsepower"].astype(int, copy=True)

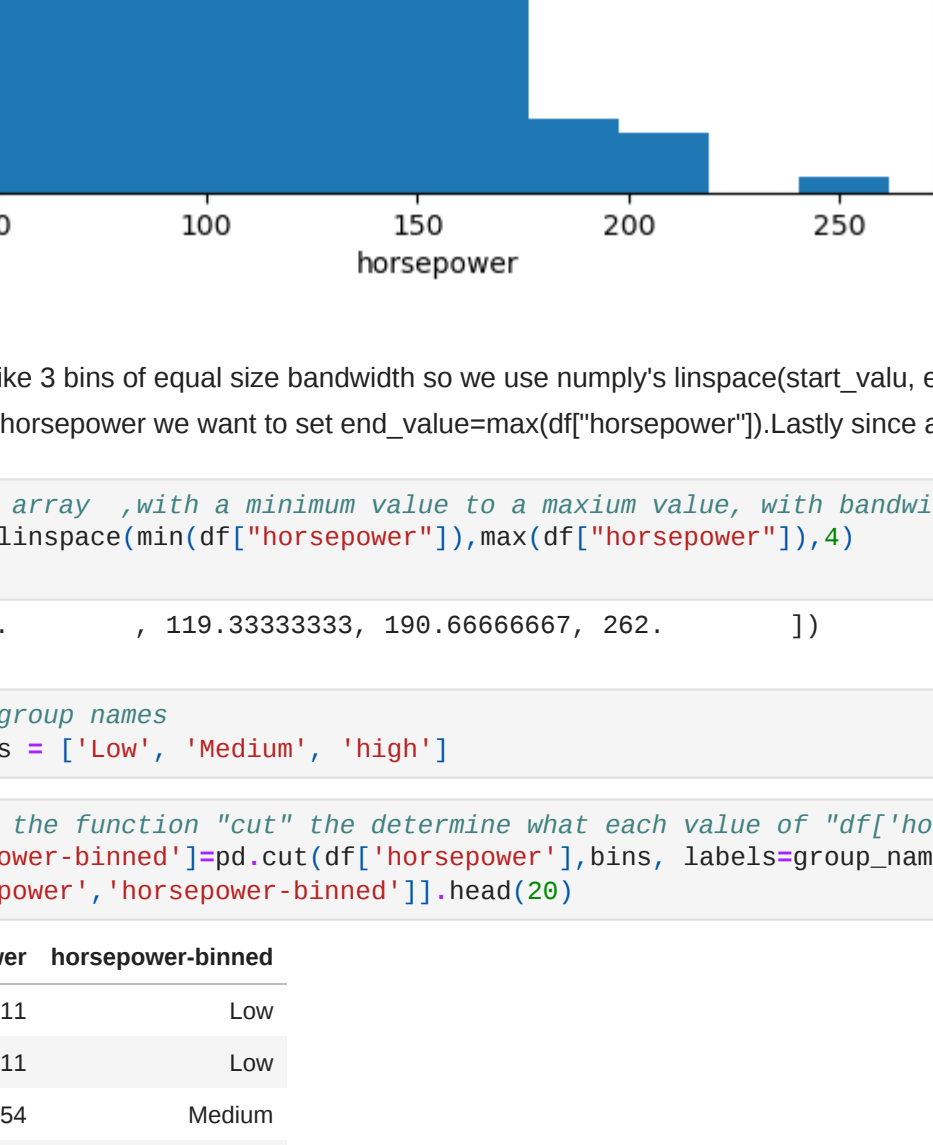
In [45]: #Histogram of horsepower ,to see what the distribution of horse looks like.
matplotlib inline
import matplotlib.pyplot as plt
from matplotlib import pyplot

plt.pyplot.hist(df["horsepower"])

# set x/y labels and plot title
plt.pyplot.xlabel("horsepower")
plt.pyplot.ylabel("count")
plt.pyplot.title("horsepower bins")

Out[45]:
Text(0.5, 1.6, 'horsepower bins')

In [46]:
horsepower bins



1. I would like 3 bins of equal size bandwidth so we use numpy's linspace(start,yaku,end_value, numbers_ generated function. Since I want to include the minimum value of horsepower we want to set start_value = min(df["horsepower"])Also i want to include the maximum value of horsepower we want to set end_value=max(df["horsepower"])Lastly since am building 3 bins of equal length ,there should be 4 dividers, so numbers_ generated =4

In [46]: # Building array ,with a minimum value to a maxime value, with bandwidth calculated above.The bins will be values used to determine when the bin ends and another begins.
bins = np.linspace(min(df["horsepower"]),max(df["horsepower"]),4)

Out[46]:
array([ 48. , 119.33333333, 199.66666667, 262. ])

In [47]: #setting group names
group_names = ['Low','Medium','high']

In [48]: # Applying the function "cut" the determine what each value of "df["horsepower"]" belongs to.
df["horsepower-binned"]pd.cut(df["horsepower"],bins, labels=group_names, include_lowest = True )
df["horsepower-binned"].head(28)

Out[48]:
horsepower horsepower-binned
0 111 Low
1 111 Low
2 154 Medium
3 102 Low
4 115 Low
5 110 Low
6 110 Low
7 101 Low
8 140 Medium
9 101 Low
10 101 Low
11 121 Medium
12 121 Medium
13 121 Medium
14 182 Medium
15 182 Medium
16 182 Medium
17 48 Low
18 70 Low
19 70 Low

In [49]: # The number of vehicles in each bin
df["horsepower-binned"].value_counts()

Out[49]:
Low 153
Medium 43
High 5
Name: horsepower-binned, dtype: int64

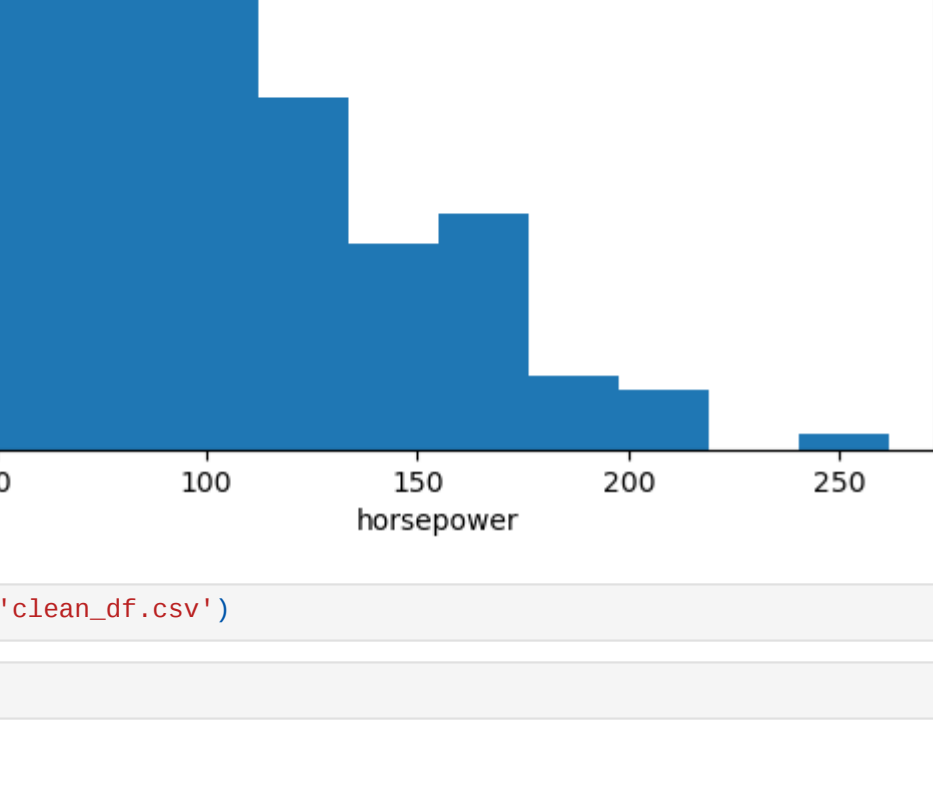
In [50]: #Plotting the distribution of each bin.
matplotlib inline
import matplotlib.pyplot as plt
from matplotlib import pyplot

plt.pyplot.hist(df["horsepower"])

# set x/y labels and plot title
plt.pyplot.xlabel("horsepower")
plt.pyplot.ylabel("count")
plt.pyplot.title("horsepower bins")

Out[50]:
Text(0.5, 1.6, 'horsepower bins')

In [51]:
horsepower bins


```