

# Performing EDA on Car information dataset



## *Importing Required Libraries*

In [1]:

```
1 # Importing Libraries
2 import pandas as pd
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5 import numpy as np
6
7 from mpl_toolkits.mplot3d import Axes3D
8
9 # Import warnings to ignore any warnings
10 import warnings
11 warnings.filterwarnings('ignore')
```

## Insert dataset

In [2]:

```
1 # Load dataset
2 data = pd.read_csv('Vehicle Manufacturing.csv')
```

`read_csv()` function : Read a comma-separated values (csv) file into DataFrame.

---

## Displaying data

In [3]:

```
1 # Display data
2 data
```

Out[3]:

		name	mpg	cylinders	displacement	horsepower	weight	acceleration	model_year	ori
0		chevrolet chevelle malibu	18.0	8	307.0	130.0	3504	12.0	70	
1		buck skylark 320	15.0	8	350.0	165.0	3693	11.5	70	
2		plymouth satellite	18.0	8	318.0	150.0	3436	11.0	70	
3		amc rebel sst	16.0	8	304.0	150.0	3433	12.0	70	
4		ford torino	17.0	8	302.0	140.0	3449	10.5	70	
...		...	...	...	...	...	...	...	...	...
393		ford mustang gl	27.0	4	140.0	86.0	2790	15.6	82	
394		vw pickup	44.0	4	97.0	52.0	2130	24.6	82	eur
395		dodge rampage	32.0	4	135.0	84.0	2295	11.6	82	
396		ford ranger	28.0	4	120.0	79.0	2625	18.6	82	
397		chevy s- 10	31.0	4	119.0	82.0	2720	19.4	82	

398 rows × 9 columns



## Displaying first 5 data

In [4]:

```
1 # Display the first 5 rows
2 data.head()
```

Out[4]:

		name	mpg	cylinders	displacement	horsepower	weight	acceleration	model_year	origin
0		chevrolet chevelle malibu	18.0	8	307.0	130.0	3504	12.0	70	usa
1		buick skylark 320	15.0	8	350.0	165.0	3693	11.5	70	usa
2		plymouth satellite	18.0	8	318.0	150.0	3436	11.0	70	usa
3		amc rebel sst	16.0	8	304.0	150.0	3433	12.0	70	usa
4		ford torino	17.0	8	302.0	140.0	3449	10.5	70	usa

- 
- To view first 5 rows we use `head()` function
  - This function returns the first `n` rows for the object based on position.
  - It is useful for quickly testing if your object has the right type of data in it.
  - If n is larger than the number of rows, this function returns all rows.
-

## Displaying last 5 rows

In [5]:

```
1 # Display the Last 5 rows
2 data.tail()
```

Out[5]:

		name	mpg	cylinders	displacement	horsepower	weight	acceleration	model_year	ori
393		ford mustang gl	27.0	4	140.0	86.0	2790	15.6	82	1
394		vw pickup	44.0	4	97.0	52.0	2130	24.6	82	euro
395		dodge rampage	32.0	4	135.0	84.0	2295	11.6	82	1
396		ford ranger	28.0	4	120.0	79.0	2625	18.6	82	1
397		chevy s-10	31.0	4	119.0	82.0	2720	19.4	82	1

- To view last 5 rows we use `tail()` function

- This function returns last `n` rows from the object based on position.

- It is useful for quickly verifying data, for example, after sorting or appending rows.

---

## Display more information about all attributes.

In [6]:

```
1 # Get more information about dataset
2 data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 398 entries, 0 to 397
Data columns (total 9 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   name            398 non-null    object 
 1   mpg             398 non-null    float64
 2   cylinders       398 non-null    int64  
 3   displacement    398 non-null    float64
 4   horsepower      392 non-null    float64
 5   weight          398 non-null    int64  
 6   acceleration    398 non-null    float64
 7   model_year      398 non-null    int64  
 8   origin          398 non-null    object 
dtypes: float64(4), int64(3), object(2)
memory usage: 28.1+ KB
```

- This method prints information about a DataFrame including the index dtype and columns, non-null values and memory usage.

### Observation

- We get columns name, non-null count, data type
  - To get more information about dataset we use info() function
- 

## Changing the datatype

In [7]:

```
1 # The following lines convert the column in the DataFrame to the 'category'
2 data['cylinders'] = data['cylinders'].astype('category')
3 data['model_year'] = data['model_year'].astype('category')
4 data['origin'] = data['origin'].astype('category')
```

### Observation

- The column is converted to the 'category' data type
- 

## To get number of rows and columns

In [8]:

```
1 # Get number of rows and columns
2 data.shape
```

Out[8]: (398, 9)

### Observation

- To get number of rows and columns we use shape function
  - number of rows are 398
  - number of columns are 9
-

## ***checking for missing values***

```
In [9]: 1 # Check for missing values  
2 data.isnull().sum()
```

```
Out[9]: name      0  
mpg       0  
cylinders 0  
displacement 0  
horsepower  6  
weight     0  
acceleration 0  
model_year  0  
origin     0  
dtype: int64
```

*Detect missing values.*

### **Observation**

- we get null values by using `isnull()` function
- There is null value in horsepower

## ***Handling null values***

```
In [10]: 1 # Handling null values  
2 data.dropna(inplace=True)
```

### **Note**

- Null Values can be handled by replacing or dropping them
- in this situation we use droping method as the values are less and cannot effect the dataset

## ***Reset index***

```
In [11]: 1 #after deleting index is set  
2 data.reset_index(drop=True, inplace=True)  
3 data.shape
```

```
Out[11]: (392, 9)
```

### **Note**

`reset_index()` method used to

*Reset the index of the DataFrame, and use the default one instead.*

DEPRECATION WARNING: This method will be removed in a future version.

In [12]:

```
1 # check for nullvalue again
2 data.isnull().sum()
```

Out[12]:

```
name          0
mpg           0
cylinders     0
displacement  0
horsepower    0
weight         0
acceleration  0
model_year    0
origin         0
dtype: int64
```

## Note

*isnull() method is use to Detect missing values.*

## Observation

- After dropping null values we can observe that all the null values are removed

---

## Generate descriptive statistics

In [13]:

```
1 # Summary statistics of numerical variables
2 data.describe()
```

Out[13]:

	mpg	displacement	horsepower	weight	acceleration
<b>count</b>	392.000000	392.000000	392.000000	392.000000	392.000000
<b>mean</b>	23.445918	194.411990	104.469388	2977.584184	15.541327
<b>std</b>	7.805007	104.644004	38.491160	849.402560	2.758864
<b>min</b>	9.000000	68.000000	46.000000	1613.000000	8.000000
<b>25%</b>	17.000000	105.000000	75.000000	2225.250000	13.775000
<b>50%</b>	22.750000	151.000000	93.500000	2803.500000	15.500000
<b>75%</b>	29.000000	275.750000	126.000000	3614.750000	17.025000
<b>max</b>	46.600000	455.000000	230.000000	5140.000000	24.800000

## Note

*describe() method Generate descriptive statistics.*

*Descriptive statistics include those that summarize the central tendency, dispersion and shape of a dataset's distribution*

*For numeric data, the result's index will include count, mean, std, min, max as well as lower, 50 and upper percentiles. By default the lower percentile is 25 and the upper percentile is 75, The*

---

## Duplicate values

In [14]:

```
1 # Check for Duplicate Value
2 data[data.duplicated()]
```

Out[14]:

```
name  mpg  cylinders  displacement  horsepower  weight  acceleration  model_year  origin
```

---

### Observation

- it check row wise data
  - there is no duplicates in dataset
- 

## Number of vehicle models year wise

In [15]:

```
1 # Display total number of vehicles
2 d = data.groupby('model_year')['name'].count().sort_values(ascending=False)
3 print(d)
```

```
model_year
73      40
78      36
76      34
75      30
82      30
70      29
79      29
72      28
77      28
81      28
71      27
80      27
74      26
Name: name, dtype: int64
```

### Observation

- The 'd' Series displays the total number of vehicles for each 'model\_year', sorted in descending order based on the count.

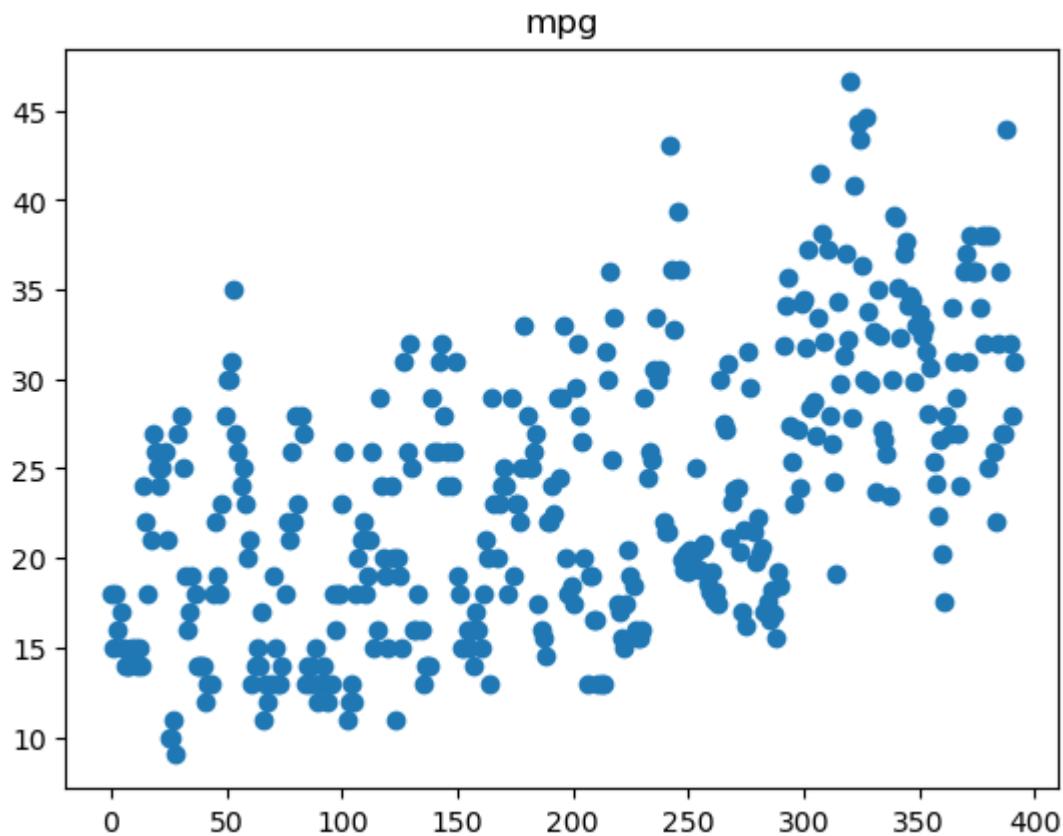
- The highest count appears at the top, indicating the year with the most vehicles, while the lowest count appears at the bottom.
- 

## Univariate Data Visualizations

### Univariate scatter plot

In [16]:

```
1 # Scatter plot of 'mpg' against the DataFrame index
2
3 plt.scatter(data.index,data[ 'mpg' ])
4 plt.title('mpg')
5 plt.show()
```



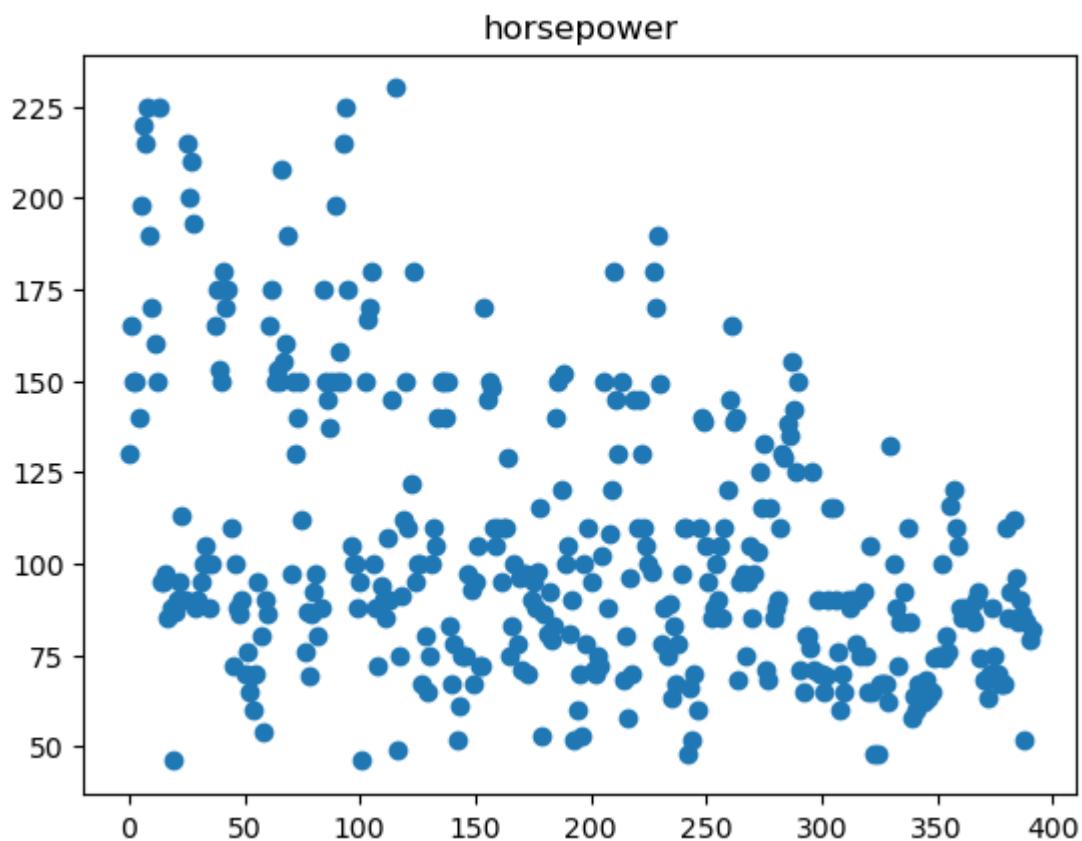
### Observation

-this plot shows the value of mpg according to there serial no.

- Scatter plots are useful for identifying patterns, trends, and potential outliers in the data.

In [17]:

```
1 # Scatter plot of 'horsepower' against the DataFrame index
2
3 plt.scatter(data.index,data[ 'horsepower' ])
4 plt.title('horsepower')
5 plt.show()
```



### Observation

- the plot shows the horse power of vehicles serial wise

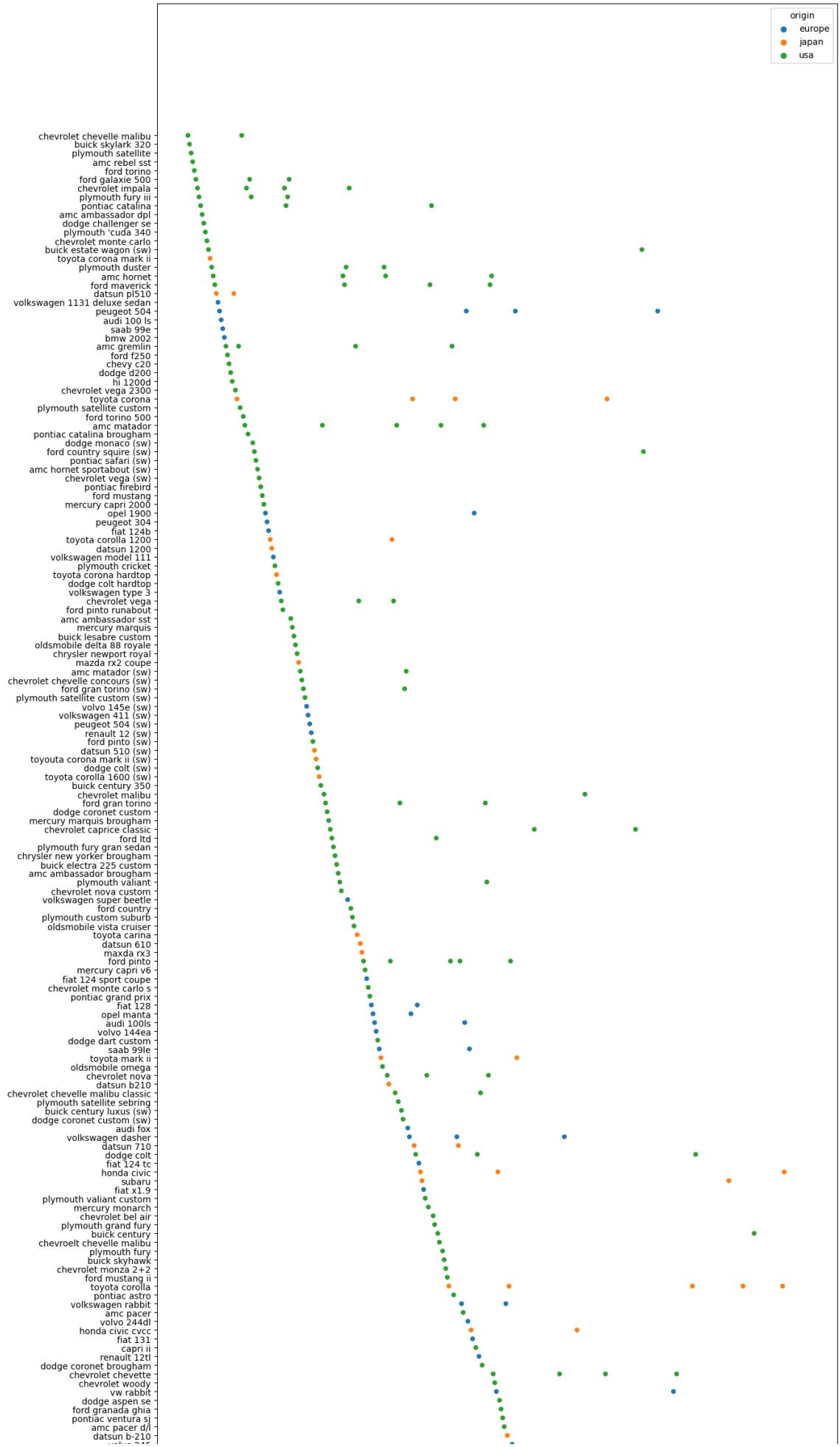
## **Univariate scatter plot by considering its origin**

In [18]:

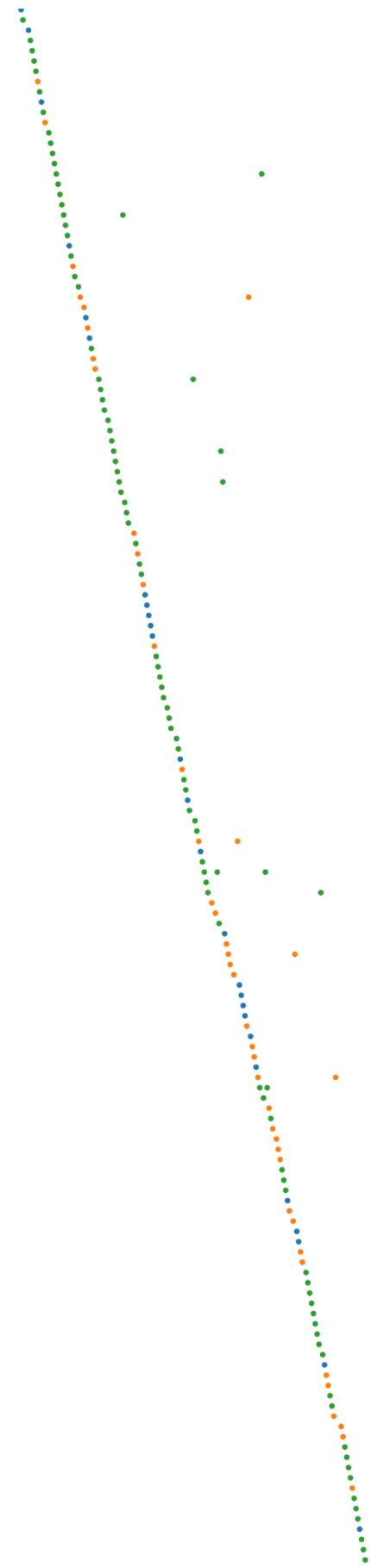
```
1 # Scatter plot of 'name' against the DataFrame index, colored by 'origin'  
2 plt.figure(figsize=(14,60))  
3 sns.scatterplot(x=data.index,y=data['name'],hue=data['origin'])
```

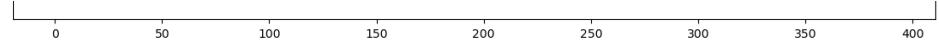
Out[18]: <Axes: ylabel='name'>





name  
volvo 242  
plymouth volare premier v8  
mercedes-benz 280s  
cadillac deville  
chevy s-10  
ford f108  
dodge d100  
honda accord cvcc  
buick opel isuzu deluxe  
renault 5 gti  
plymouth arena  
datsun f-10 hatchback  
oldsmobile cutlass supreme  
dodge monaco brougham  
mercury cougar brougham  
chevrolet concorde  
chevrolet caprice  
plymouth valiant  
ford granada  
pontiac grand prix lj  
chevrolet monte carlo landau  
chrysler cordoba  
ford thunderbird  
volkswagen rabbit  
volkswagen sunbird coupe  
toyota corolla liftback  
ford mustang ii 2+2  
dodge colt m/m  
subaru dl  
dodge omni  
bmw 320i  
mazda rx-4  
volkswagen rabbit custom diesel  
ford fiesta  
mazda glc deluxe  
datsun b210 gx  
oldsmobile cutlass  
dodge diplomat  
mercury monarch ghia  
pontiac phoenix lj  
ford fairmont (auto)  
ford fairmont (man)  
plymouth amc concord  
buick century special  
mercury zephyr  
dodge aspen  
amc concord d/l  
buick regal sport coupe (turbo)  
dodge futura  
dodge magnum xe  
datsun 510  
dodge omni  
toyota celica gt liftback  
plymouth sapporo  
oldsmobile starfire x  
datsun 200-sx  
audi 5000  
volvo 264gl  
saab 99gl  
peugeot 604sl  
volkswagen accord  
honda accord lx  
pontiac lemans v6  
mercury zephyr 6  
ford fairmont 4  
amc concord dl  
dodge aspen 6  
ford ltd landau  
mercury landau  
dodge st. regis  
chevrolet malibu classic (sw)  
chrysler lebaron town @ country (sw)  
vw rabbit custom  
maxda 626  
dodge colt hatchback  
amc spirit dl  
mercedes benz 300d  
cadillac eldorado  
plymouth horizon  
plymouth horizon tc3  
fiat strada custom  
buick skylark limited  
oldsmobile omega brougham  
pointiac phoenix  
toyota corolla tercel  
dodge omni  
ford fairmont  
audi 4000  
toyota corona liftback  
mazda 626  
datsun 510 hatchback  
toyota corolla  
vw rabbit c (diesel)  
vw dasher (diesel)  
audi 5000s (diesel)  
mercedes-benz 240d  
honda civic 1500 gl  
volkswagen rabbit  
datsun 200-sx  
mazda rx-7 gs  
triumph tr7 coupe  
honda accord  
plymouth reliant  
dodge aries wagon (sw)  
honda accord (sw)  
plymouth chano  
honda civic 1300  
datsun 210 mpg  
toyota tercel  
mazda glc 4  
plymouth colizo 4  
ford escort 4w  
ford escort 2h  
volkswagen jetta  
honda prelude  
datsun 200sx  
peugeot 505s turbo diesel  
toyota corolla  
toyota cressida  
datsun 810 maxima  
oldsmobile cutlass ls  
ford granada gl  
chrysler lebaron salon  
chevrolet cavalier  
chevrolet cavalier (pony)  
chevrolet cavalier 2-door  
pontiac j2000 se hatchback  
dodge aries se  
ford fairmont futura  
volkswagen rabbit  
mazda 626  
mercury glc custom  
plymouth horizon miser  
mercury lynx l  
nissan stanza xe  
honda civic (auto)  
datsun 110 gx  
buick century limited  
oldsmobile cutlass ciera (diesel)  
chrysler lebaron medallion  
ford granada l  
toyota celica gt  
dodge charger 2.2  
chevrolet camaro  
ford mustang gl  
vw pickup  
dodge rampage  
ford ranger  
chevy s-10





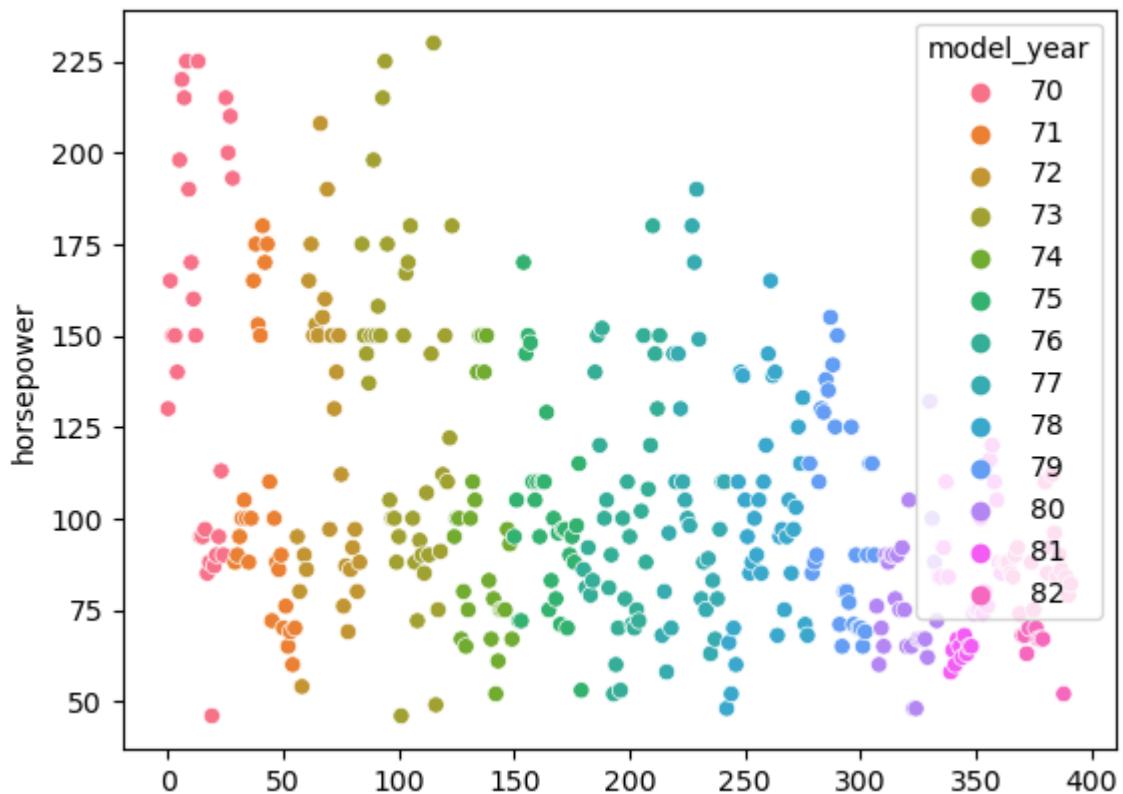
## Observation

- the colors represent where the vehicles are origin and it seems that maximum cars are made in usa

## Univariate scatter plot by considering its year

```
In [19]: 1 # Scatter plot of 'horsepower' against the DataFrame index, colored by 'model_year'
2 sns.scatterplot(x=data.index,y=data['horsepower'],hue=data['model_year'])
```

```
Out[19]: <Axes: ylabel='horsepower'>
```



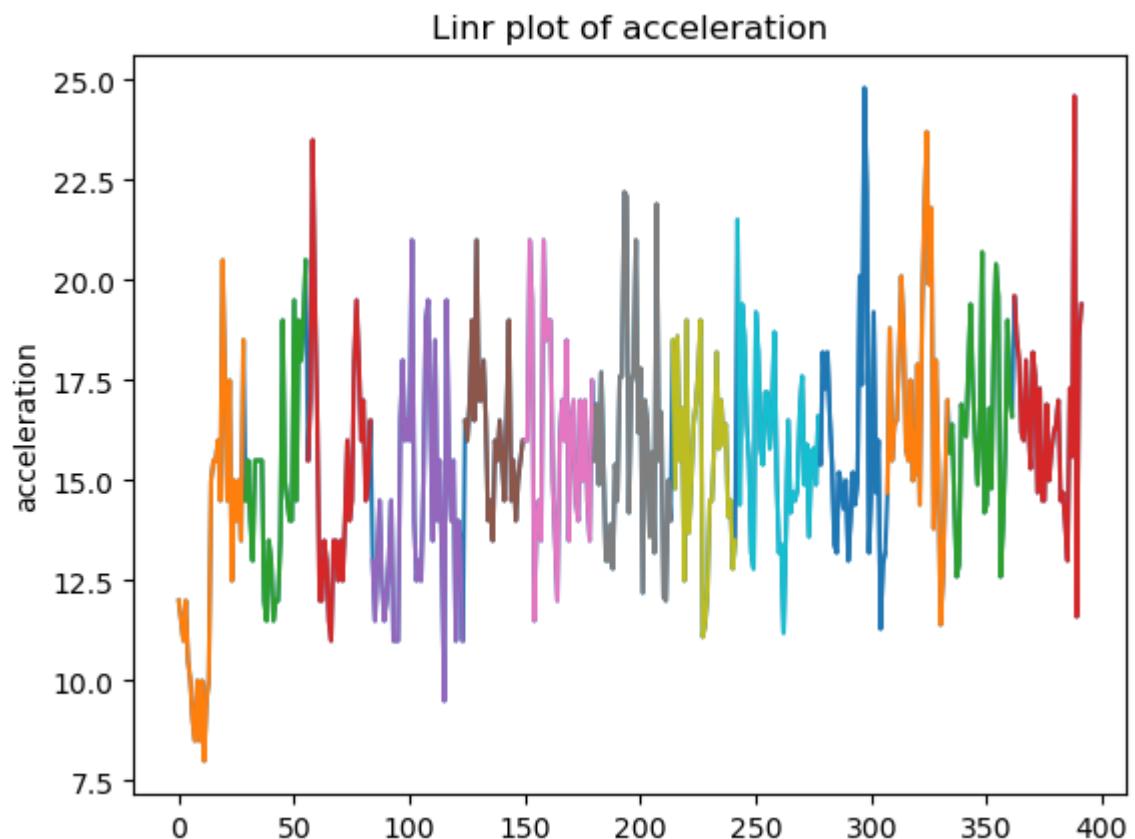
## Observation

- The scatter plot visualizes the relationship between 'horsepower' and the DataFrame index, the colors shows the the modal year

## Univariate linr plot

In [20]:

```
1 # Line plot of 'acceleration' against the DataFrame index, grouped by 'model_year'
2
3 plt.title('Linr plot of acceleration')
4 plt.ylabel('acceleration')
5
6 # Plotting the Line plot for the entire dataset
7 plt.plot(data.index,data['acceleration'])
8
9 # Grouping the data by 'model_year' and plotting Line plots for each group
10 for name,group in data.groupby('model_year'):
11     plt.plot(group.index,group['acceleration'])
12
13 plt.show()
```



### Observation

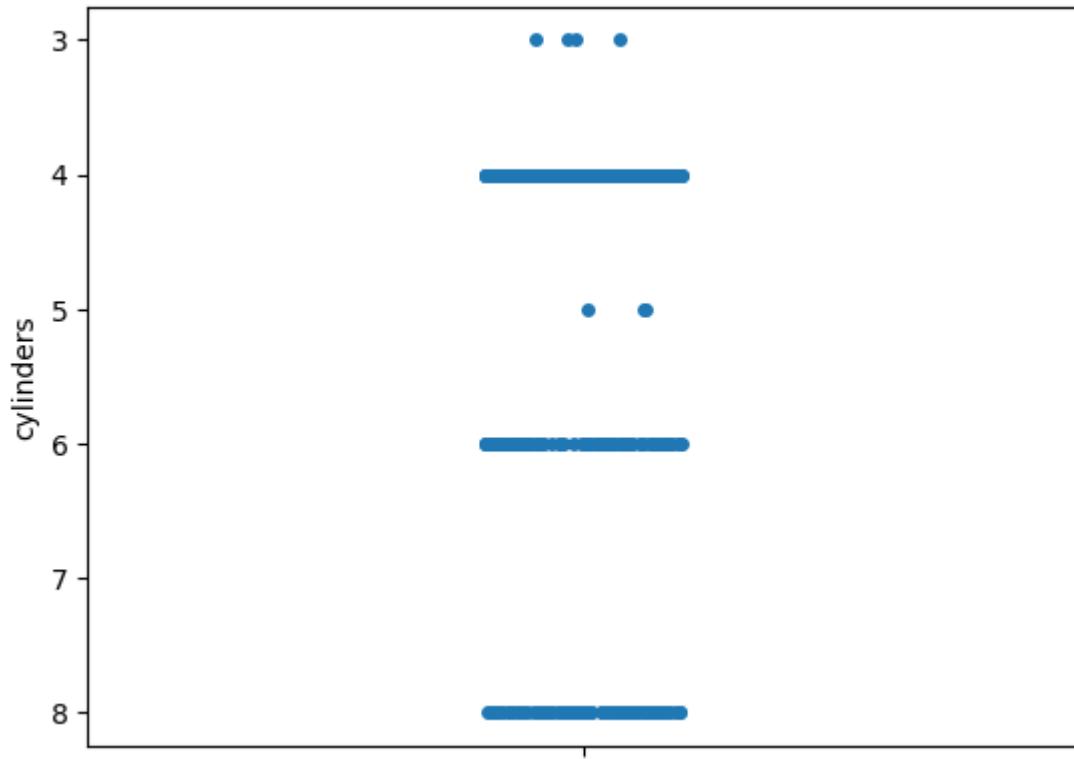
- The code above creates a line plot of the 'acceleration' values against the DataFrame index.

## Univariate Strip plot

In [21]:

```
1 # Strip plot of 'cylinders'  
2 sns.stripplot(y=data['cylinders'])
```

Out[21]: <Axes: ylabel='cylinders'>

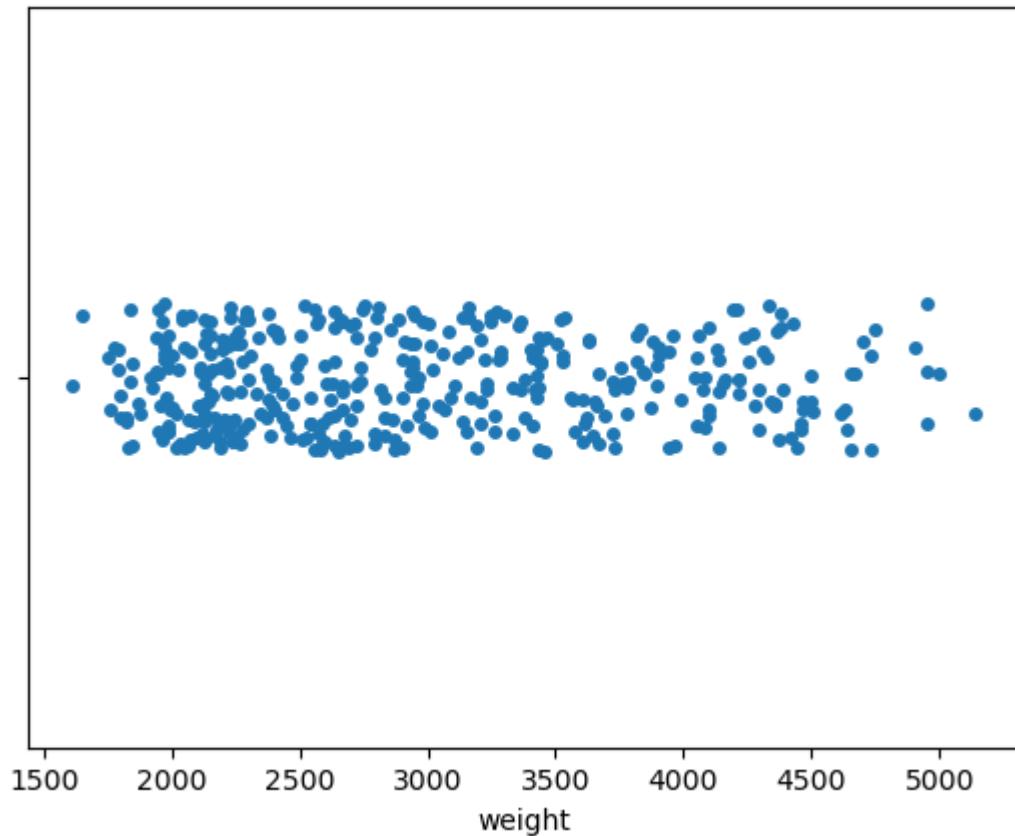


### Observation

- The strip plot visualizes the distribution of the 'cylinders' variable, which likely represents the number of cylinders in vehicles, as it is plotted on the y-axis.

```
In [22]: 1 # Strip plot of 'weight'  
2 sns.stripplot(x=data['weight'])
```

```
Out[22]: <Axes: xlabel='weight'>
```



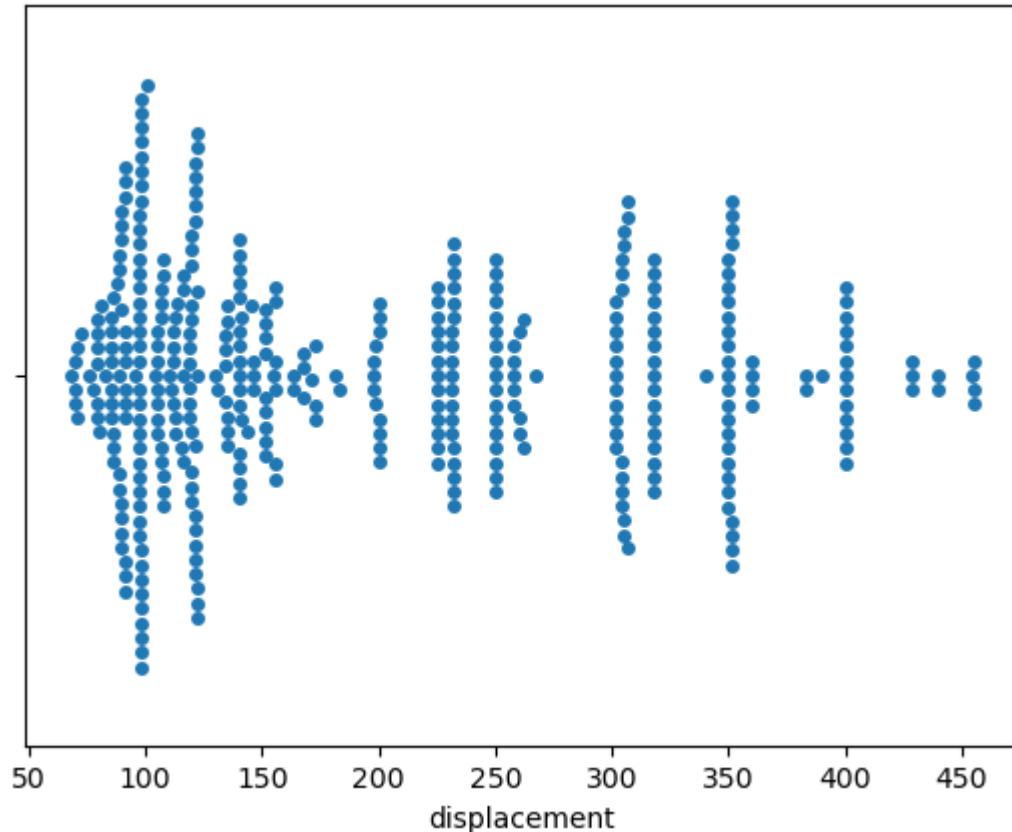
### Observation

- The scatter plot visualizes the relationship between 'weight' and the DataFrame index, allowing us to examine how the 'weight' values are distributed along the index.

## Univariate swarm plot

```
In [23]: 1 # performing swamp plot of 'displacement'  
2 sns.swarmplot(x= data['displacement'])
```

Out[23]: <Axes: xlabel='displacement'>

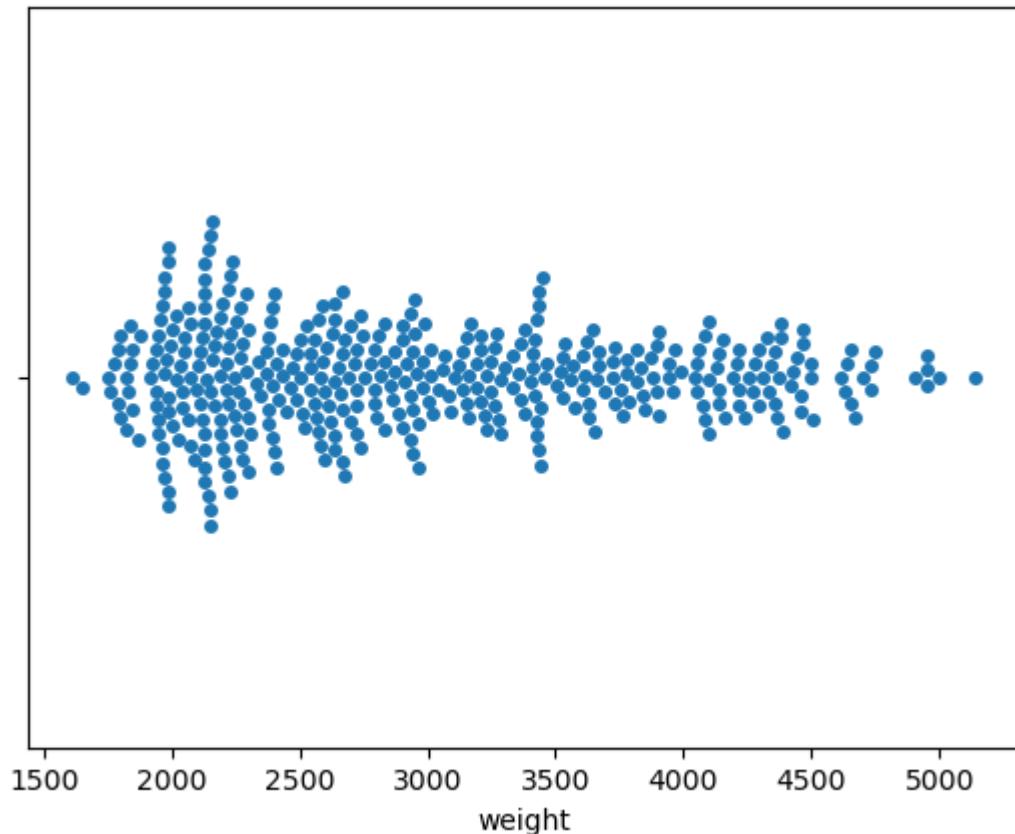


### Observation

- The code above creates a swarm plot of the 'displacement' values. The 'displacement' values are plotted on the x-axis.

```
In [24]: 1 # performing swamp plot of 'weight'  
2 sns.swarmplot(x=data['weight'])
```

Out[24]: <Axes: xlabel='weight'>



### Observation

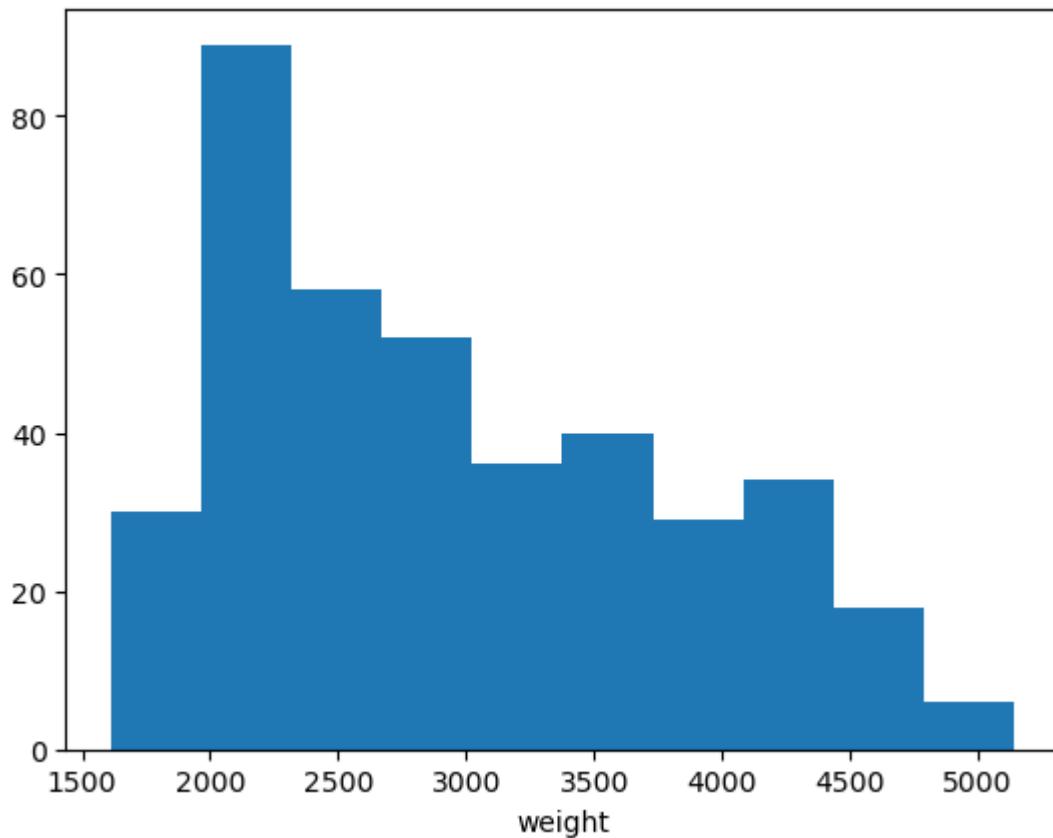
- The code above creates a swarm plot of the 'weight' values. The 'weight' values are plotted on the x-axis.
- The swarm plot is a categorical scatter plot that displays individual data points for the 'weight' variable along a single axis.

## **Univariate Histogram()**

In [25]:

```
1 # Performing histogram of 'weight'  
2 plt.hist(x=data['weight'])  
3 plt.xlabel('weight')
```

Out[25]: Text(0.5, 0, 'weight')

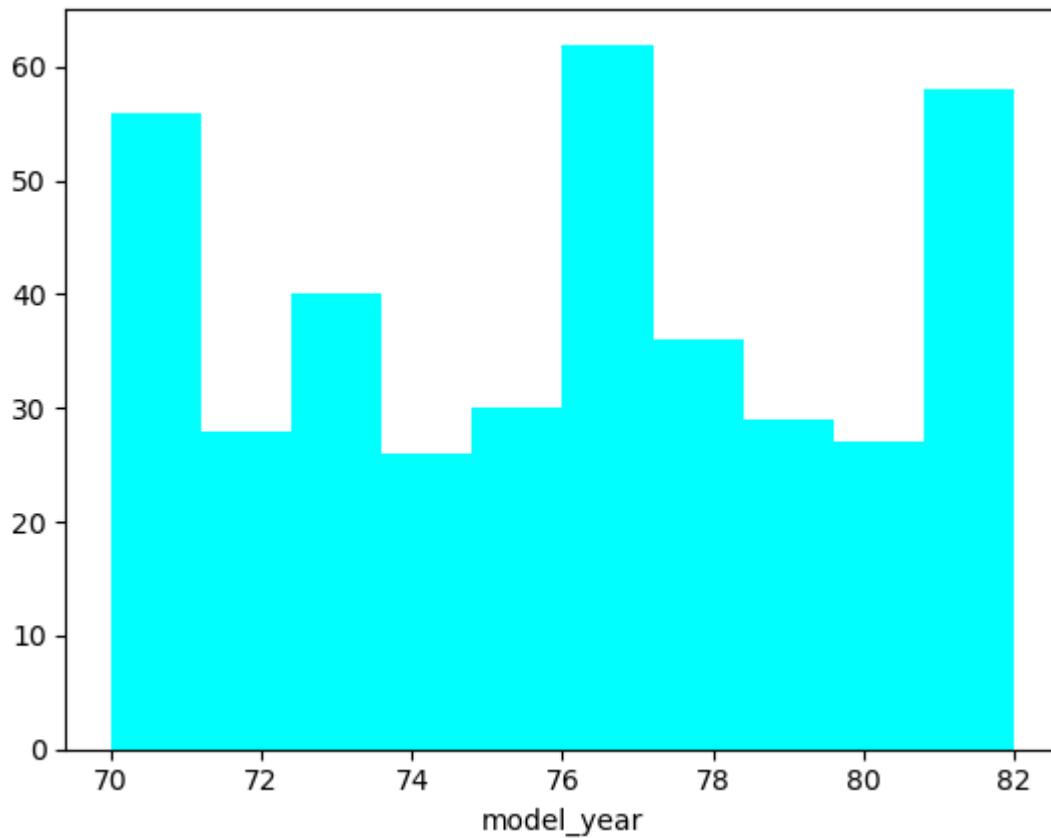


### **Observation**

- The histogram effectively displays the distribution of 'weight' values, showing the frequency of data points falling within different 'weight' bins or intervals.
- Histograms are useful for visualizing the underlying distribution of continuous variables, such as 'weight', and provide insights into the data's central tendency and spread.

```
In [26]: 1 # Histogram of 'model_year'  
2 plt.hist(x=data['model_year'],color='cyan')  
3 plt.xlabel('model_year')
```

```
Out[26]: Text(0.5, 0, 'model_year')
```



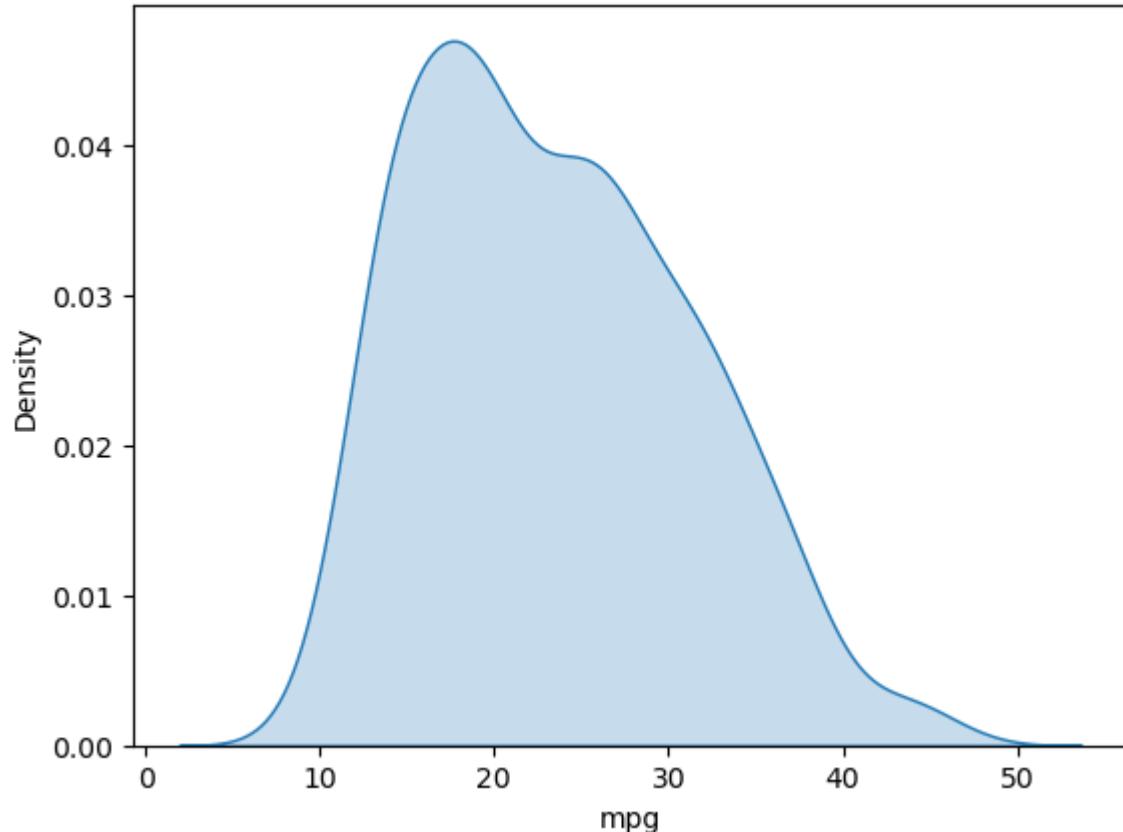
### Observation

- The histogram effectively displays the distribution of 'model\_year' values, showing the frequency of data points for each year.
- The x-axis represents the 'model\_year' values, and the y-axis represents the frequency or count of data points corresponding to each year.

## Univariate Density plot

```
In [27]: 1 # Kernel Density Estimation (KDE) plot of 'mpg'  
2 sns.kdeplot(data['mpg'], shade=True)
```

```
Out[27]: <Axes: xlabel='mpg', ylabel='Density'>
```



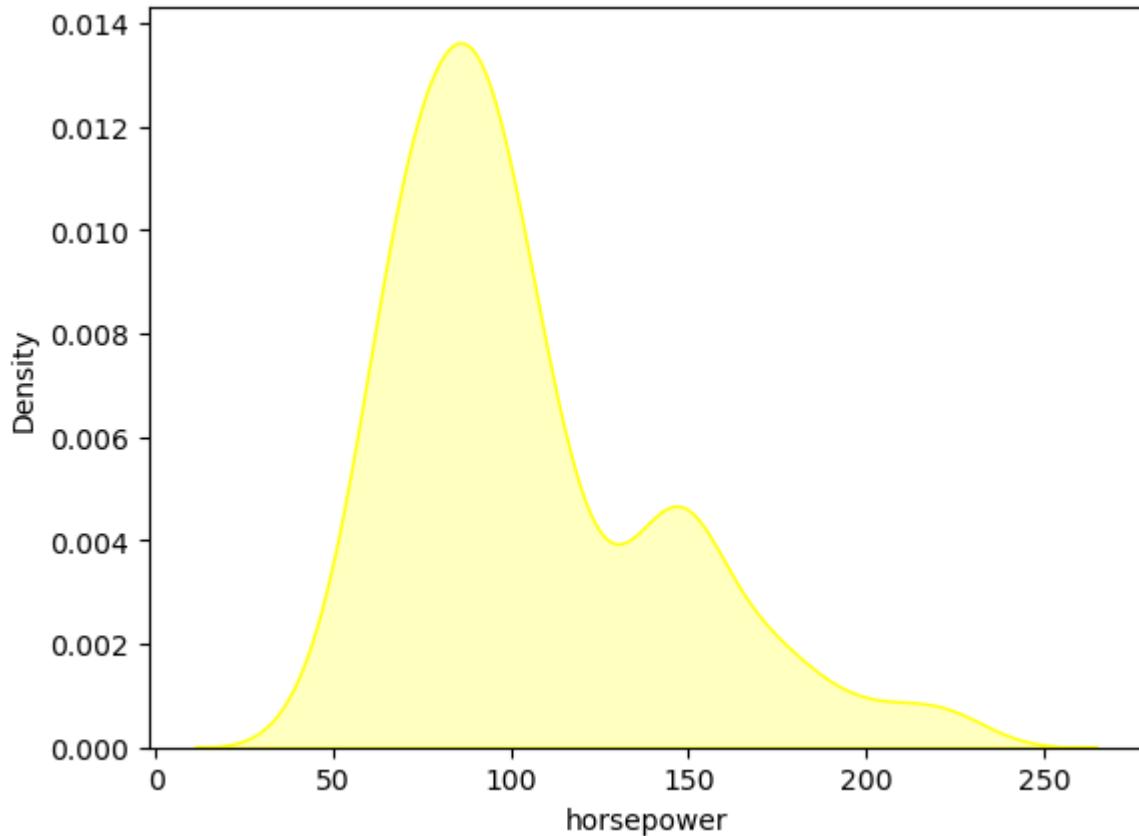
### Observation

- The KDE plot provides a smooth and continuous estimation of the probability density function of the 'mpg' variable.
- The x-axis represents the 'mpg' values, and the y-axis represents the estimated density of data points at each 'mpg' value.

In [28]:

```
1 # Kernel Density Estimation (KDE) plot of 'horsepower'  
2 sns.kdeplot(data['horsepower'], shade=True, color='yellow')
```

Out[28]: <Axes: xlabel='horsepower', ylabel='Density'>



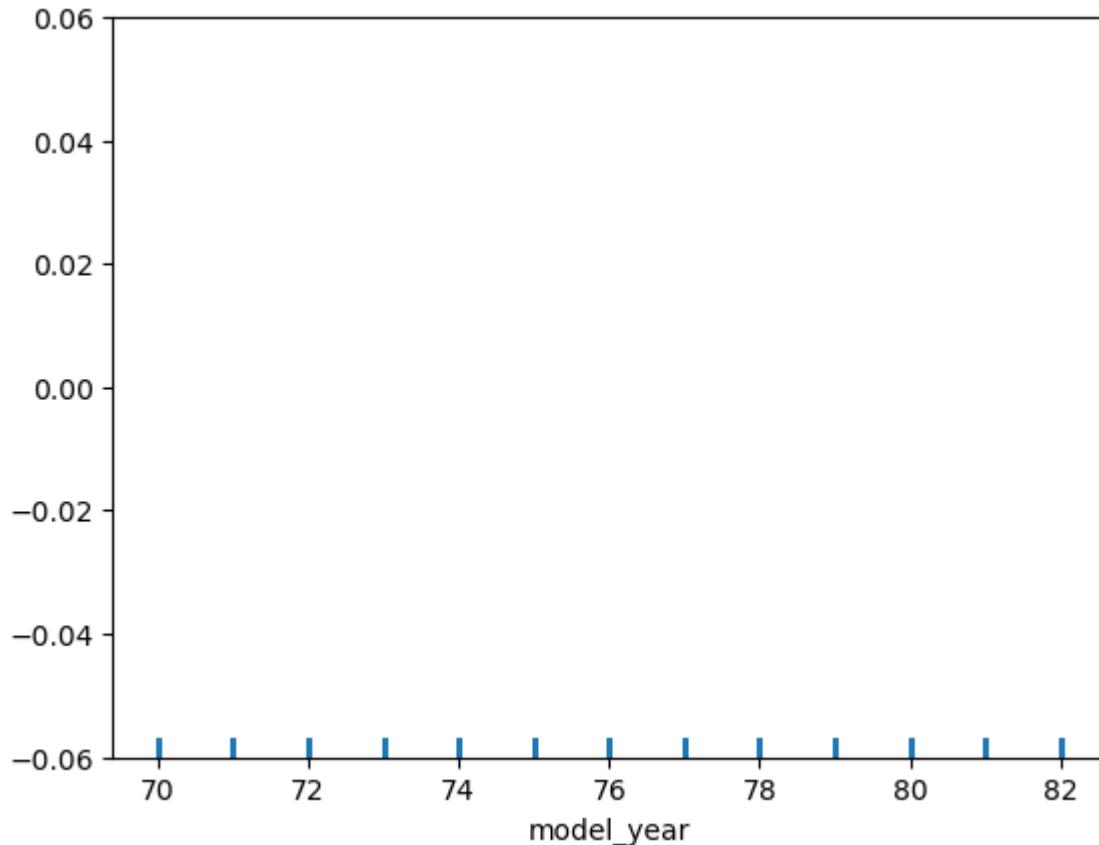
### Observation

- The 'shade=True' parameter fills the area under the curve, making the plot visually appealing and easier to interpret. The 'color='yellow"' parameter sets the color of the KDE curve to yellow.
- KDE plots are useful for visualizing the underlying distribution of continuous variables, such as 'horsepower', without the need for specifying bin sizes (unlike histograms).

## Univariate RUG plot

```
In [29]: 1 # Rug plot of 'model_year'  
2 sns.rugplot(data['model_year'])
```

```
Out[29]: <Axes: xlabel='model_year'>
```

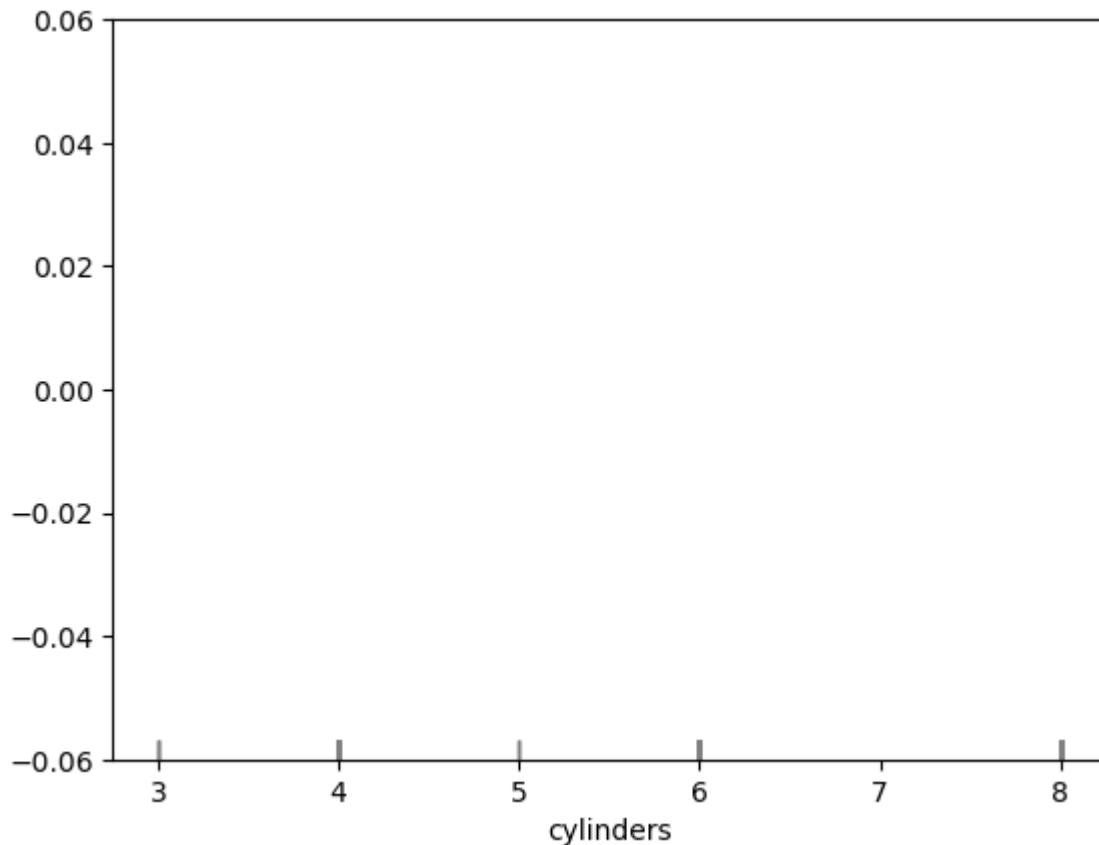


### Observation

- The rug plot is a simple one-dimensional plot that displays individual data points as ticks along a single axis.
- The x-axis represents the 'model\_year' values, and each tick (or "rug") along the axis indicates the occurrence of a single data point at that specific year.

```
In [30]: 1 # Rug plot of 'model_year'  
2 sns.rugplot(data['cylinders'],color='grey')
```

```
Out[30]: <Axes: xlabel='cylinders'>
```



### Observation

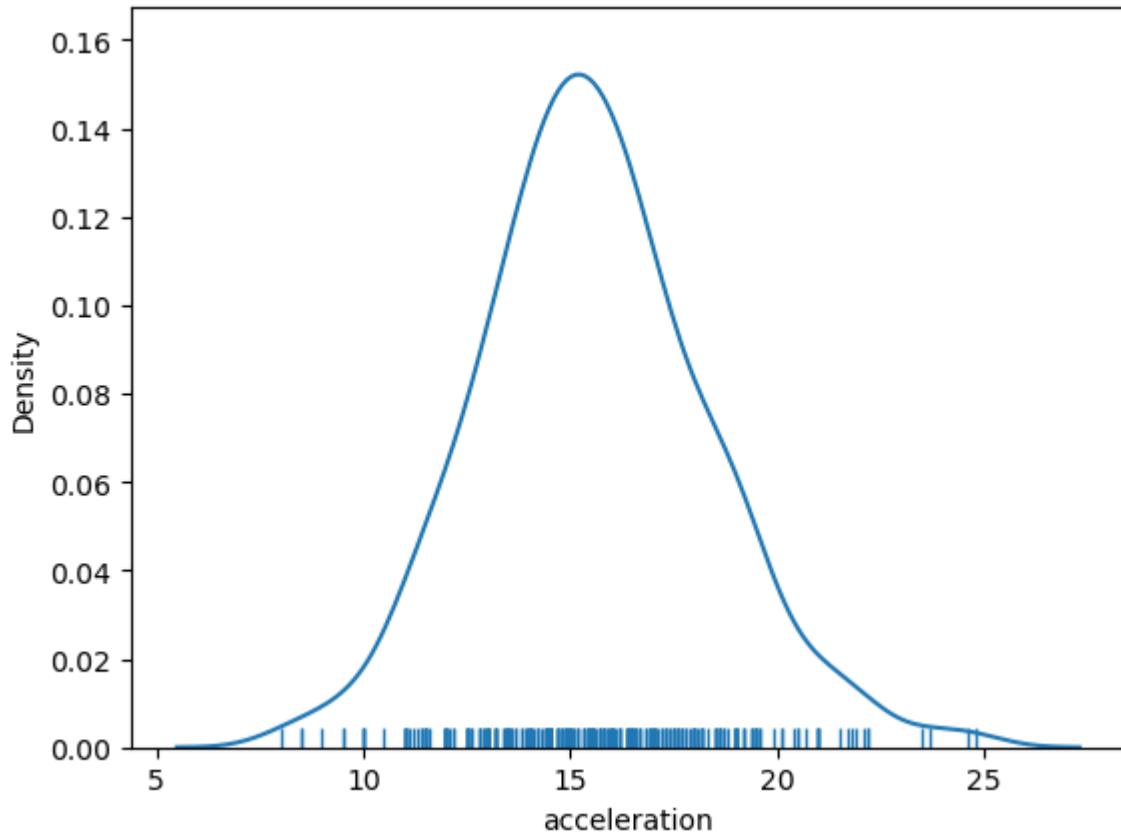
- Rug plots are used to visualize the distribution and density of data points along an axis, providing a basic overview of the data's positions.
- The rug plot is useful for understanding the distribution of values and quickly identifying the presence and density of data points for each year.

## Univariate Kernel Density plot along with rugs

In [31]:

```
1 # density plot for 'acceleration' with rug plot
2 sns.distplot(data['acceleration'],rug=True,hist=False)
```

Out[31]: <Axes: xlabel='acceleration', ylabel='Density'>



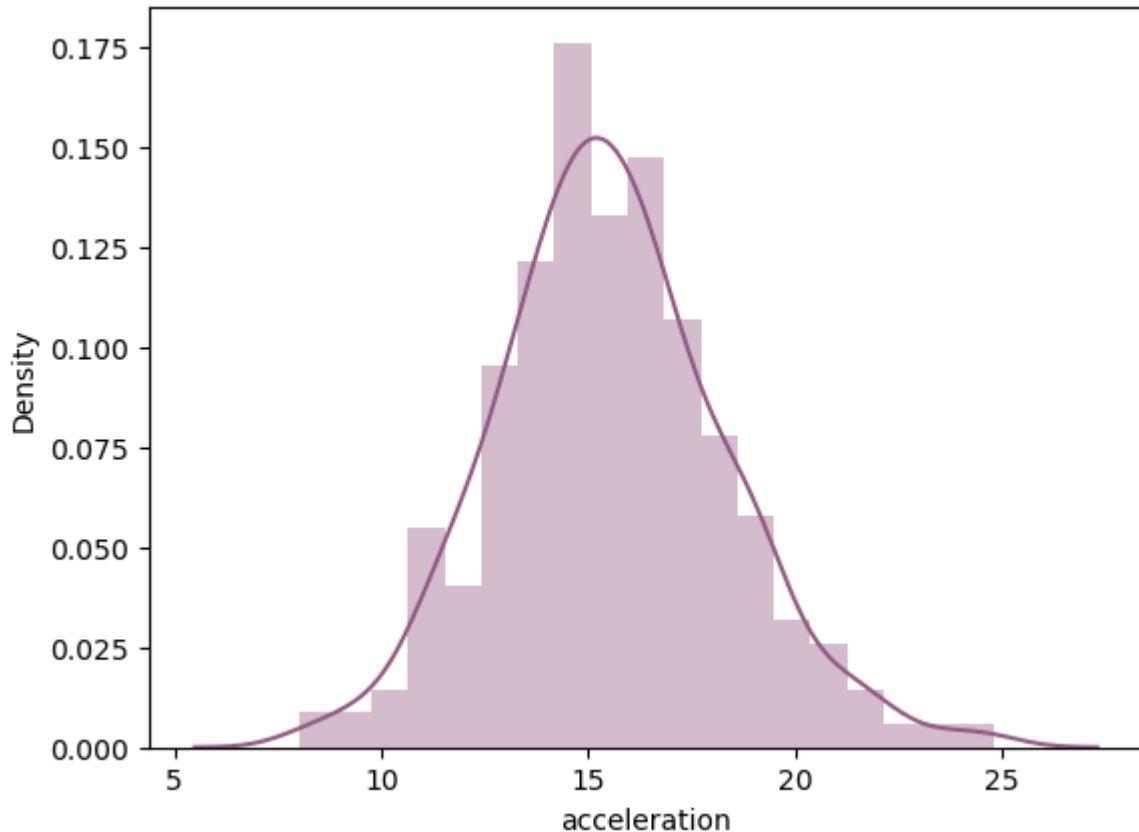
### Observation

- The distribution plot, combined with the rug plot, provides a visual representation of the 'acceleration' data distribution on a continuous axis.
- The rug plot displays individual data points as ticks along the x-axis, indicating the occurrence of each 'acceleration' value.

## Univariate Distplot

```
In [32]: 1 # Distribution plot of 'acceleration' with a custom color
2 sns.distplot(data['acceleration'],color="#915c83")
```

```
Out[32]: <Axes: xlabel='acceleration', ylabel='Density'>
```



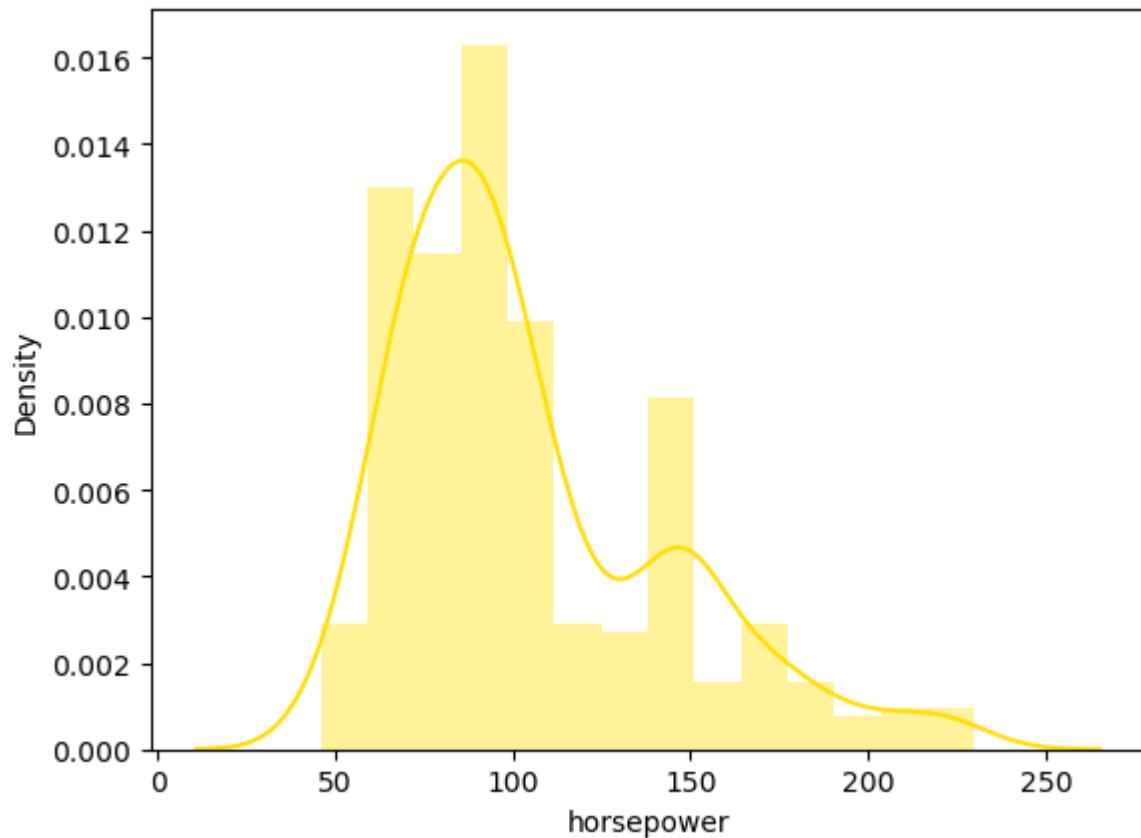
### Observation

- The distribution plot effectively displays the distribution of 'acceleration' values using a combination of a histogram, kernel density estimation (KDE) curve, and rug plot.
- The histogram bars show the binned frequency of 'acceleration' values, providing insights into the data's central tendency and spread.
- The KDE curve provides a smooth estimation of the probability density function, offering a continuous representation of the 'acceleration' distribution.

In [33]:

```
1 # Distribution plot of 'horsepower' with a custom color
2 sns.distplot(data['horsepower'],color="#ffd900")
```

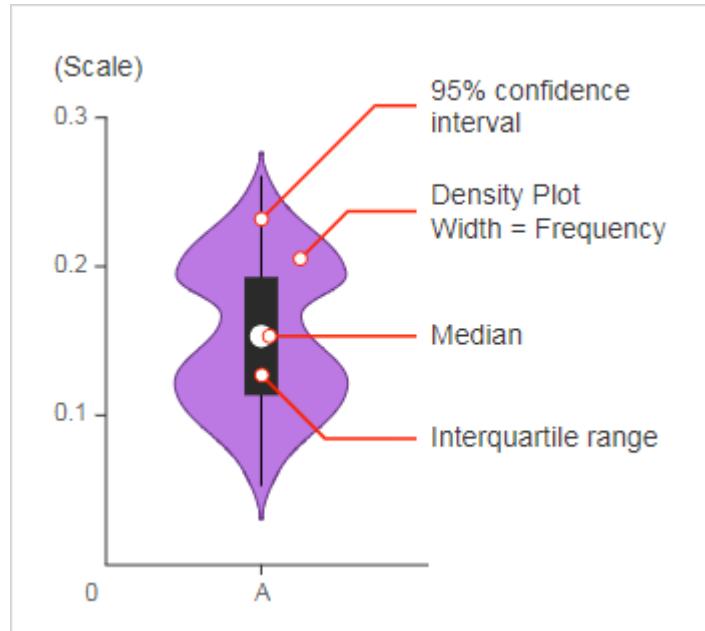
Out[33]: <Axes: xlabel='horsepower', ylabel='Density'>



### Observation

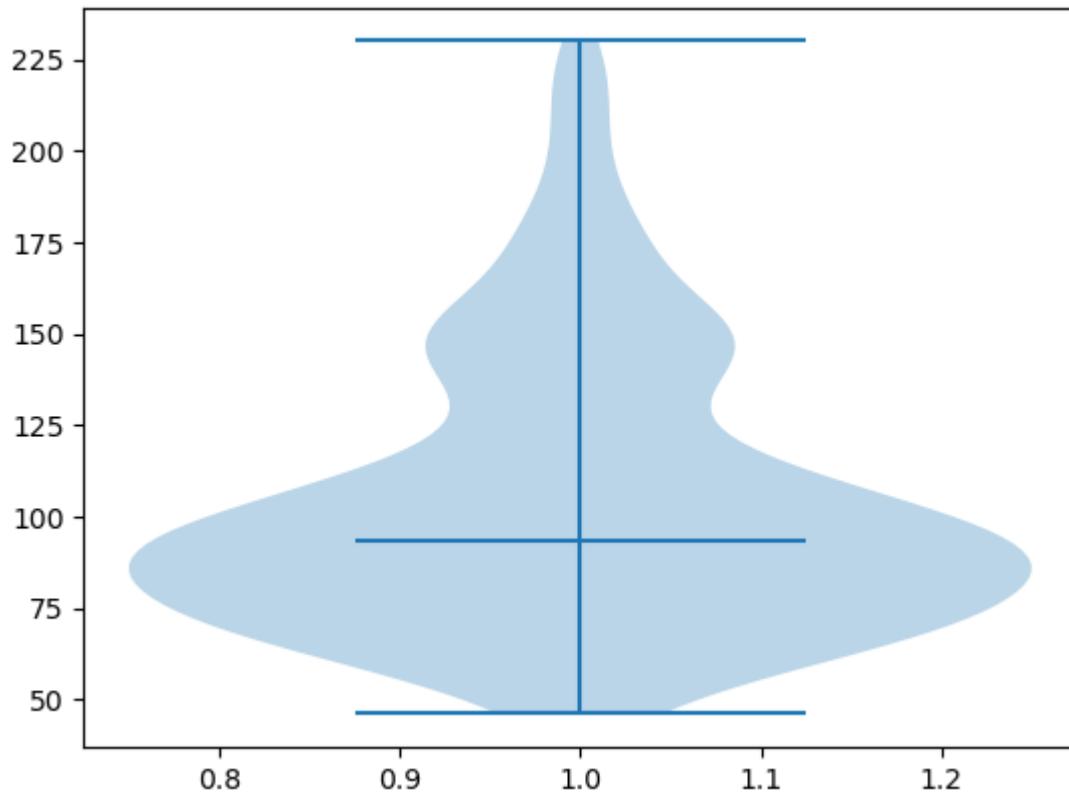
- The distribution plot with the custom color '#ffd900' is a visually appealing and informative visualization to explore the distribution of 'horsepower' values in the dataset.
- is an essential step in the exploratory data analysis process.

## **Univariate Violin Plots**



```
In [34]: 1 # Code: Violin plot of 'horsepower' with median markers displayed
          2 plt.violinplot(data.horsepower, showmedians=True)
```

```
Out[34]: {'bodies': [],
           'cmaxes': [],
           'cmins': [],
           'cbars': [],
           'cmedians': []}
```

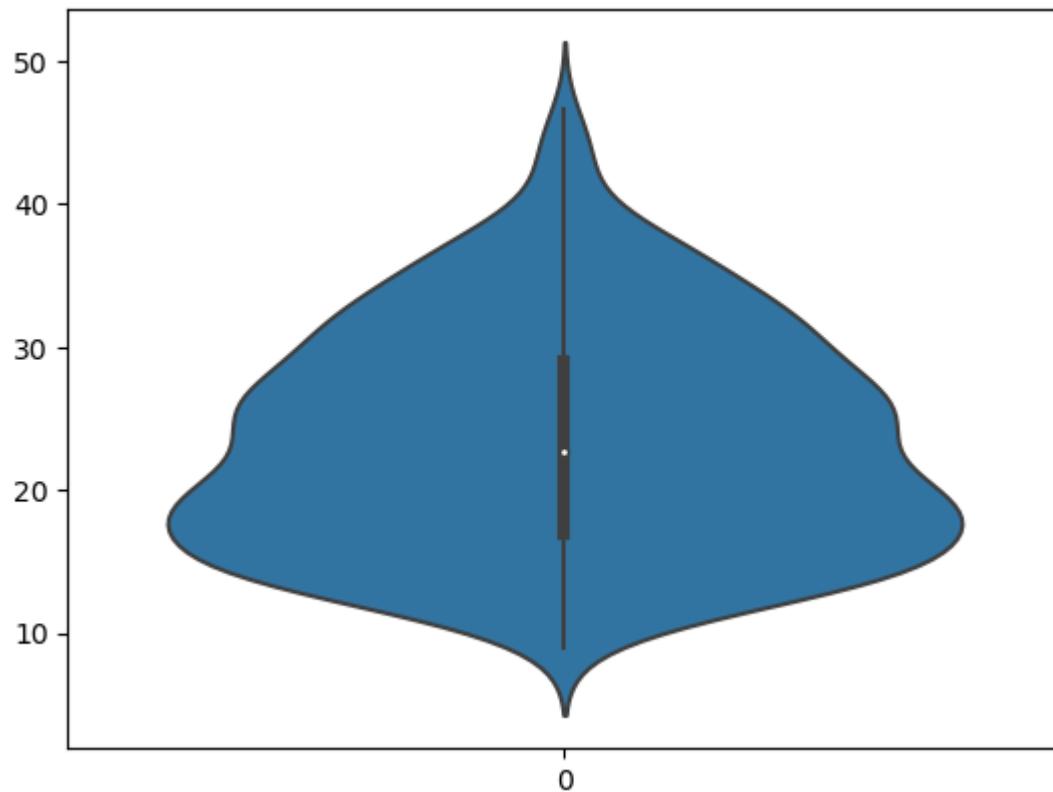


### Observation

- It combines the benefits of a box plot and a kernel density plot
- The median markers displayed inside the violins further help understand the central tendency of the data.

```
In [35]: 1 # Violin plot of 'mpg'  
2 sns.violinplot(data['mpg'])
```

Out[35]: <Axes: >



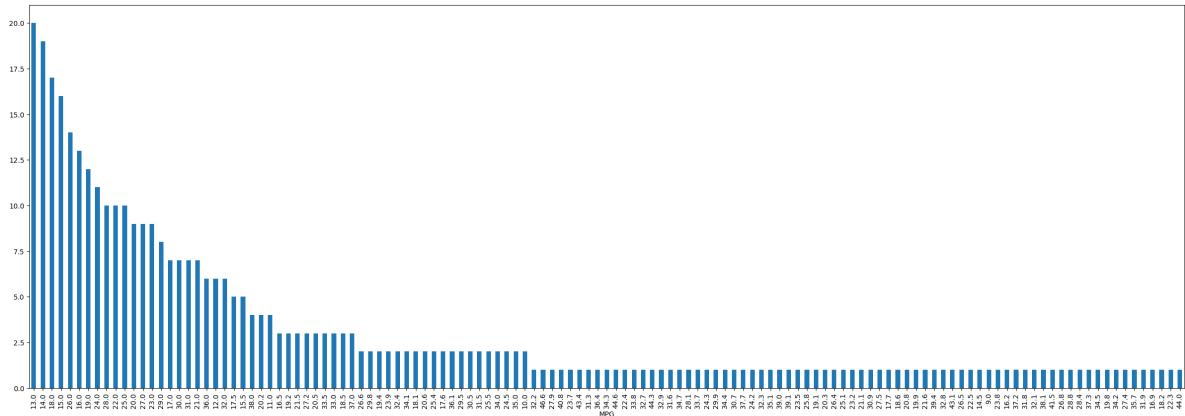
### Observation

- The plot is particularly effective for understanding the density and variability of 'mpg' across different data points.

## Univariate Bar Plot

```
In [36]: 1 # Bar plot of 'mpg' value counts
2
3 plt.figure(figsize=(30,10))
4 plt.xlabel('MPG')
5 data['mpg'].value_counts().plot.bar()
```

```
Out[36]: <Axes: xlabel='MPG'>
```



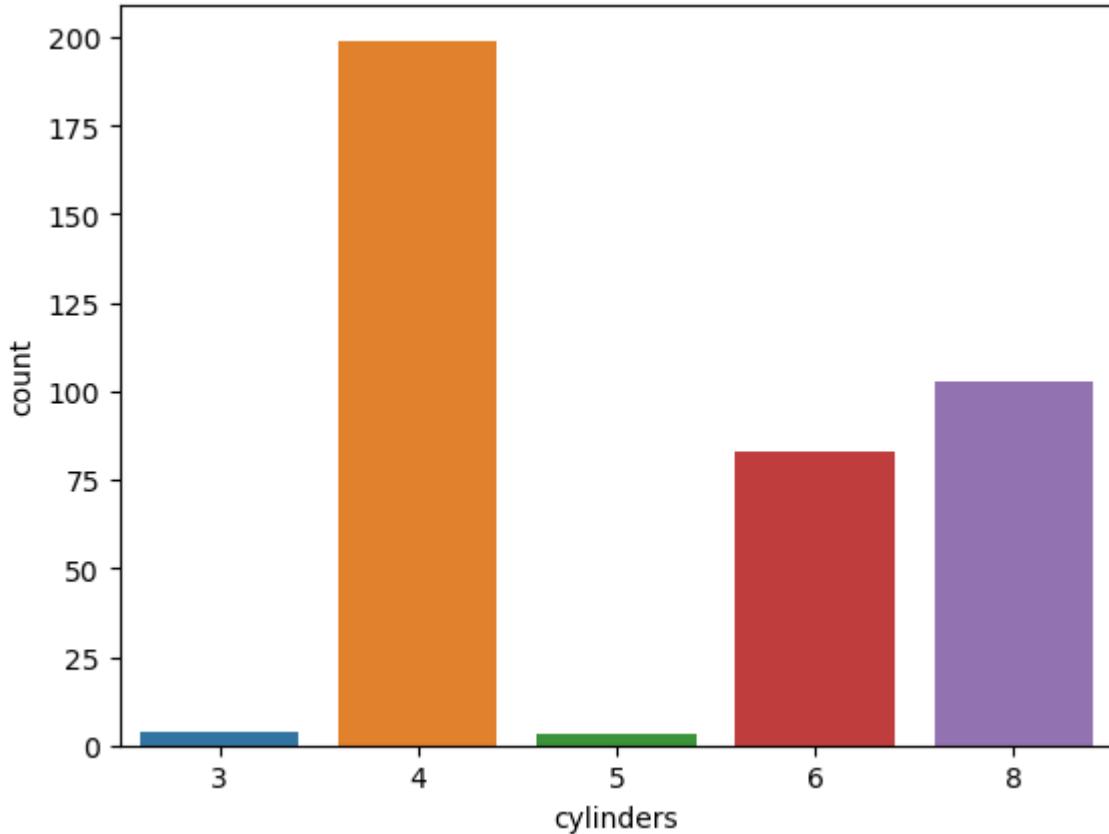
### Observation

- The bar plot displays the count of each unique 'mpg' value as vertical bars along the x-axis.

In [37]:

```
1 # Count plot of 'cylinders' using Seaborn
2 sns.countplot(x=data['cylinders'])
3 print(data['cylinders'].value_counts())
4 print((data['cylinders'].value_counts()/398)*100)
```

```
4    199
8    103
6     83
3      4
5      3
Name: cylinders, dtype: int64
4    50.000000
8    25.879397
6    20.854271
3    1.005025
5    0.753769
Name: cylinders, dtype: float64
```



### Observation

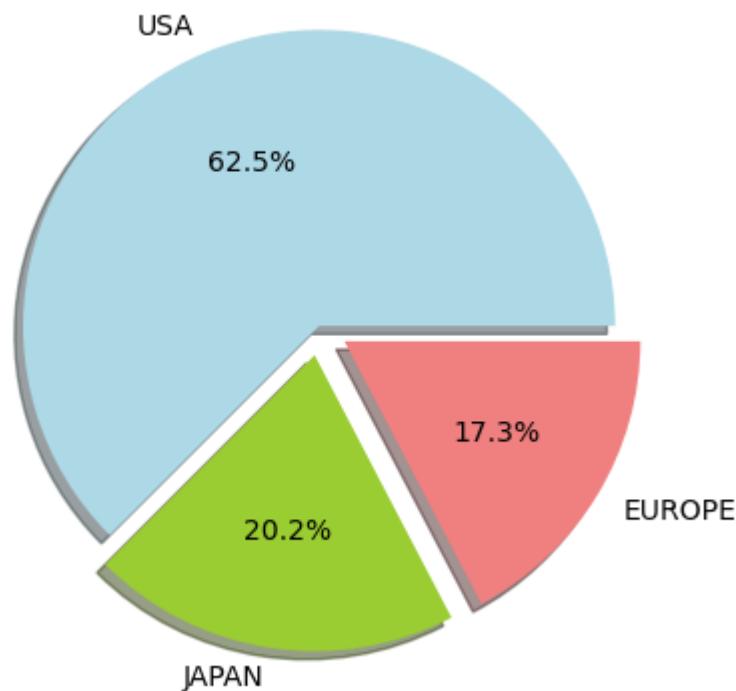
- The count plot helps us understand the distribution and prevalence of each category in the dataset.

From the count values and percentages, we can observe the relative frequency of each category in the dataset.

## **Univariate Pie chart**

In [38]:

```
1 # pie chart of origin
2 # data
3 lables=[ 'USA' , 'JAPAN' , 'EUROPE' ]
4 x=[245,79,68]
5 colors=['lightblue','yellowgreen','lightcoral','red','lightcoral']
6 explode=(0,0.1,0.1)
7
8 #plot
9 plt.pie(x,explode=explode,labels=lables,colors=colors,autopct='%1.1f%%',s
10 plt.show()
```



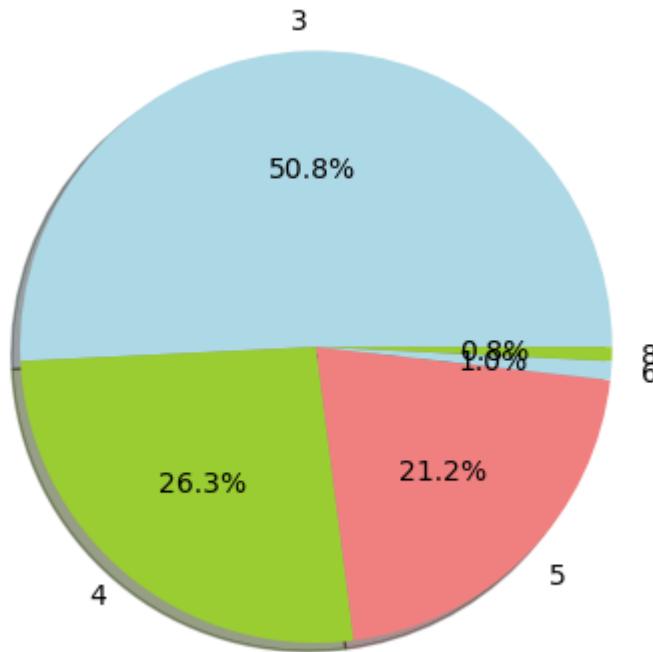
### **Observation**

- A pie chart is the most common way used to visualize the numerical proportion occupied by each of the categories.

in this pie chart we can observe number of vehicles origin in which country

In [39]:

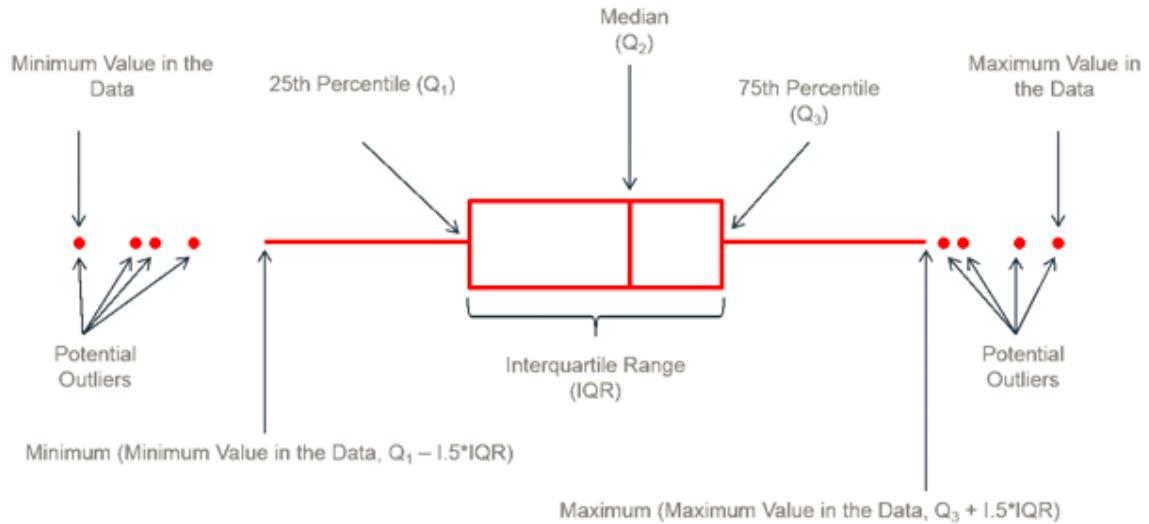
```
1 # pie chart of number of cylinders
2 # data
3 labels=[ '3', '4', '5', '6', '8']
4 x=[199,103,83,4,3]
5 colors=['lightblue','yellowgreen','lightcoral']
6 explode=(0,0,0,0,0)
7
8 #plot
9 plt.pie(x=explode,labels=labels,colors=colors,autopct='%1.1f%%',startangle=90)
10 plt.show()
```



### Observation

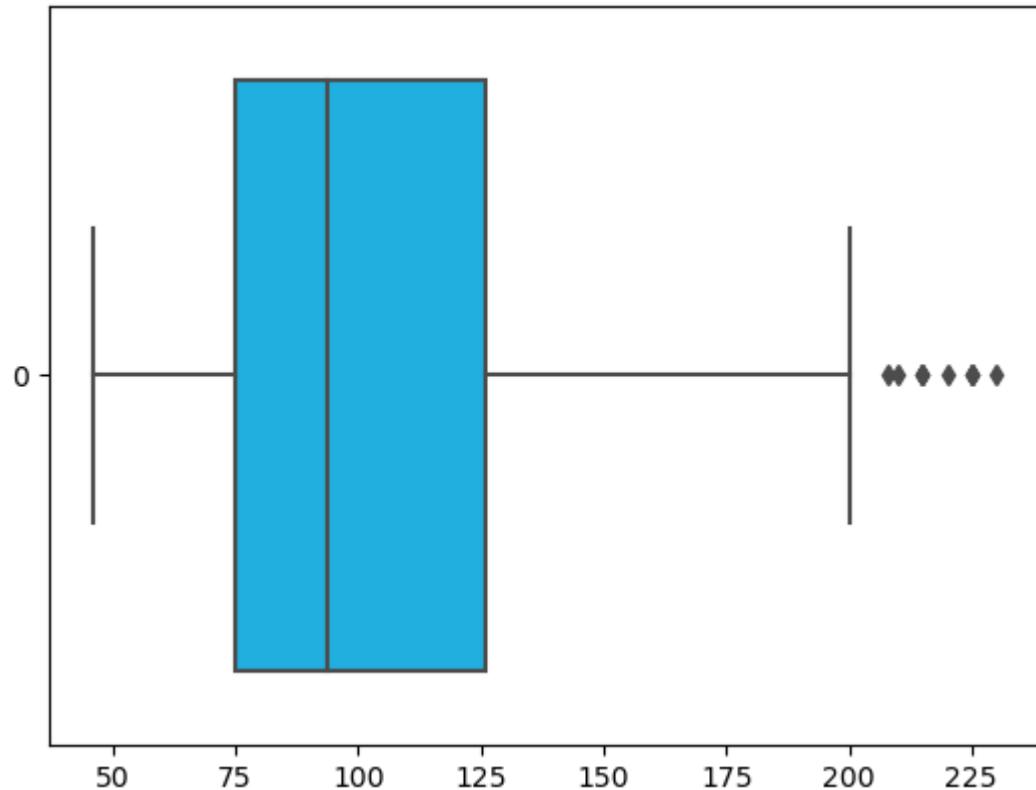
- Use the `plt.pie()` function to plot a pie chart.
- in this pie chart we can observe number of cylinders in vehicles

## Univariate Box Plot



```
In [40]: 1 # Box plot of horsepower  
2 sns.boxplot(data['horsepower'], orient='horizontal', color='deepskyblue')
```

Out[40]: <Axes: >



### Observation

- The box plot displays the distribution of the 'horsepower' variable in the dataset.

- The box represents the interquartile range (IQR), and the line inside the box represents the median value of the data

The whiskers extend from the box to the minimum and maximum non-outlier data points.

Any data points outside the whiskers are considered potential outliers and plotted as individual

### Handling the outliers using flooring and capping method

In [41]:

```
1 # flooring and capping method
2
3 outlier_list = []
4
5 # sorting data
6 data['horsepower'] = sorted(data['horsepower'])
7
8 # Finding first and third quartile
9 q1 = np.percentile(data['horsepower'], 25)
10 q3 = np.percentile(data['horsepower'], 75)
11 print("The Val of Q1 and Q2",q1, q3)
12
13 IQR = q3-q1
14
15 lwr_bound = q1-(1.5*IQR)
16 upr_bound = q3+(1.5*IQR)
17
18 # print("The Lower & Upper Bound",lwr_bound, upr_bound)
19
20 for i in data['horsepower']:
21     if (i<lwr_bound or i>upr_bound):
22         data['horsepower'] = data['horsepower'].replace(i,np.median(data[
```

The Val of Q1 and Q2 75.0 126.0

### Observation

- The provided code is implementing the flooring and capping method to handle outliers in the 'horsepower' column of the data.
- The flooring and capping method is a simple approach to handle outliers by replacing them with the median value

```
In [42]: 1 # verifying outlier
2
3 # Initialize a flag to check if there are any outliers
4 outlier_present = False
5
6 for i in data['horsepower']:
7     if (i<lwr_bound or i>upr_bound):
8         print(i)
9         outlier_present = True
10
11 # If the flag remains False, no outliers were found
12 if not outlier_present:
13     print("There are no outliers.")
```

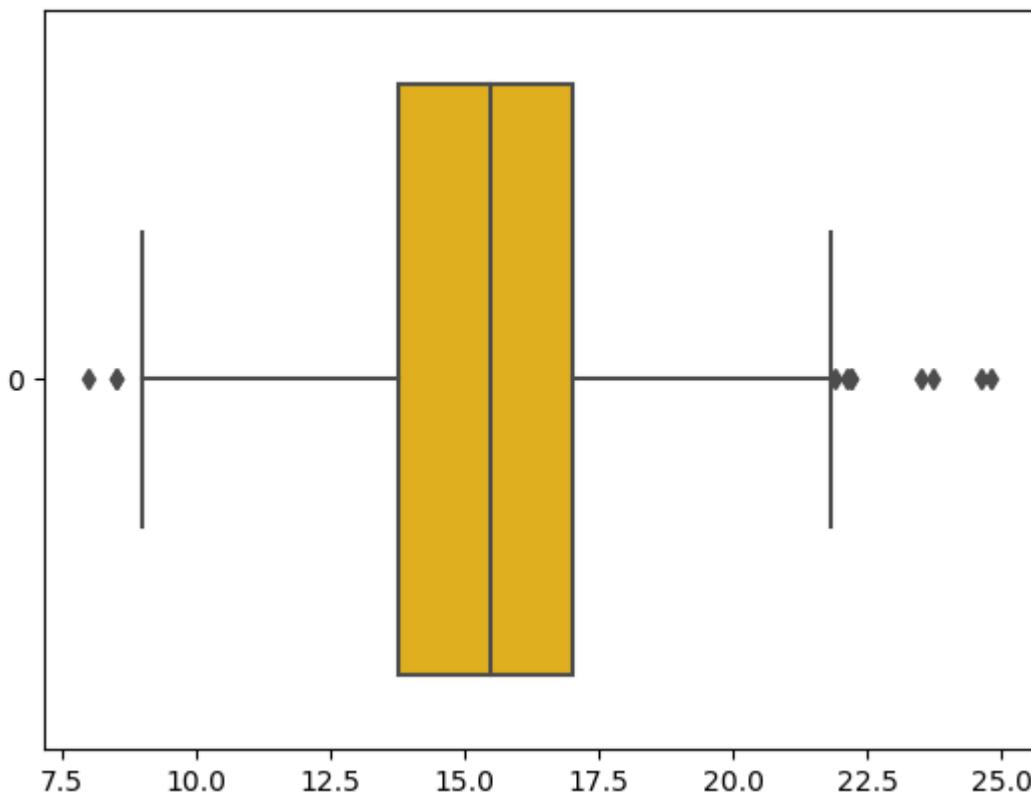
There are no outliers.

### Observation

- The provided code is verifying the presence of outliers in the 'horsepower' column of the data using the lower bound (*lwr\_bound*) and upper bound (*upr\_bound*) values calculated previously with the flooring and capping method.
- If it finds any outliers, it prints those outlier values. If no outliers are found, it informs that there are no outliers in the 'horsepower' column.

```
In [43]: 1 # box plot of acceleration
2 sns.boxplot(data['acceleration'],orient='horizontal',color='#ffbf00')
```

Out[43]: <Axes: >



## Observation

- The provided code generates a horizontal box plot for the 'acceleration' variable using Seaborn's boxplot function.

## Handling the outliers using Trimming and Capping method

```
In [44]: 1 # Handling the outliers
2 acc = data['acceleration']
3
4 tenth_percentile = np.percentile(acc, 10)
5 ninetieth_percentile = np.percentile(acc, 90)
6
7 print('10th,tenth_percentile',tenth_percentile)
8 print('90th,tenth_percentile',ninetieth_percentile)
9
10 b = np.where(acc<tenth_percentile, tenth_percentile, acc)
11
12 b1 = np.where(b>ninetieth_percentile, ninetieth_percentile, b)
13
14 data['acceleration'] = b1
```

```
10th,tenth_percentile 12.0
90th,tenth_percentile 19.0
```

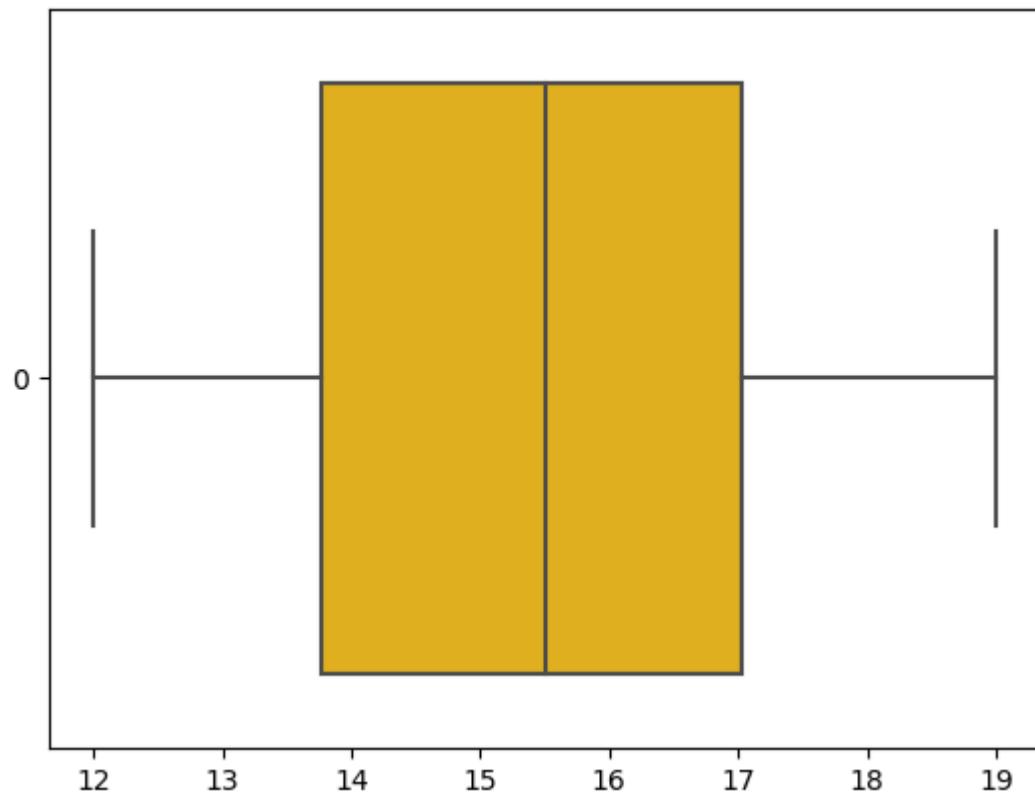
## Observation

- The provided code is handling the outliers in the 'acceleration' column of the data using a capping method

By capping the outliers within the 10th and 90th percentiles, the code ensures that extreme values are brought closer to the rest of the data

```
In [45]: 1 # verifying outliers  
2 sns.boxplot(data['acceleration'],orient='horizontal',color='#ffbfb00')
```

Out[45]: <Axes: >



### Observation

- after performing capping method, to verify it works we again run the code to perform box plot

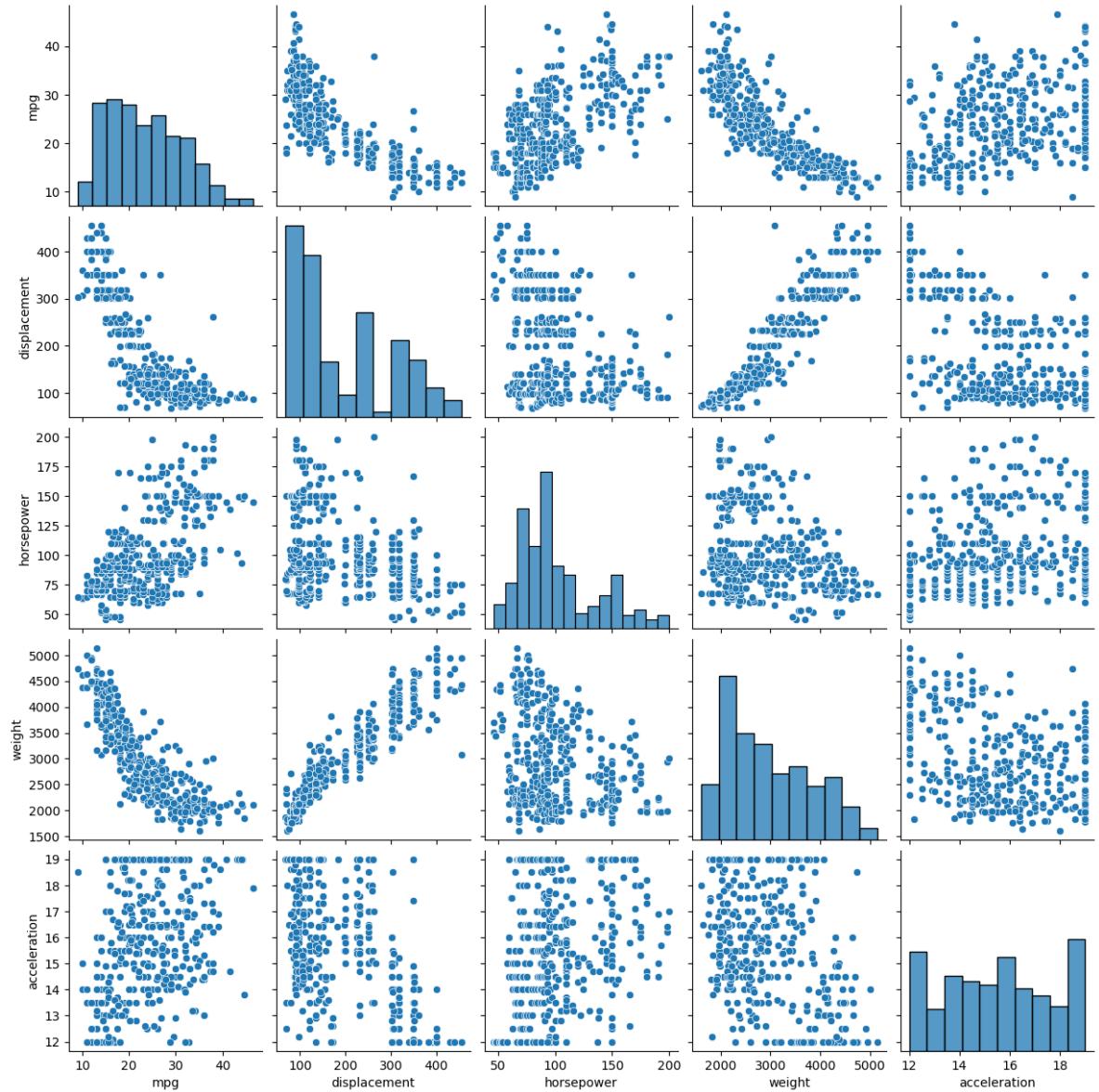
---

# Multivariate data visualization

## Scatter Plot Matrix (Pair Plot):

In [46]:

```
1 # pairplot of given dataset
2 sns.pairplot(data)
3 plt.show()
```



## Observation

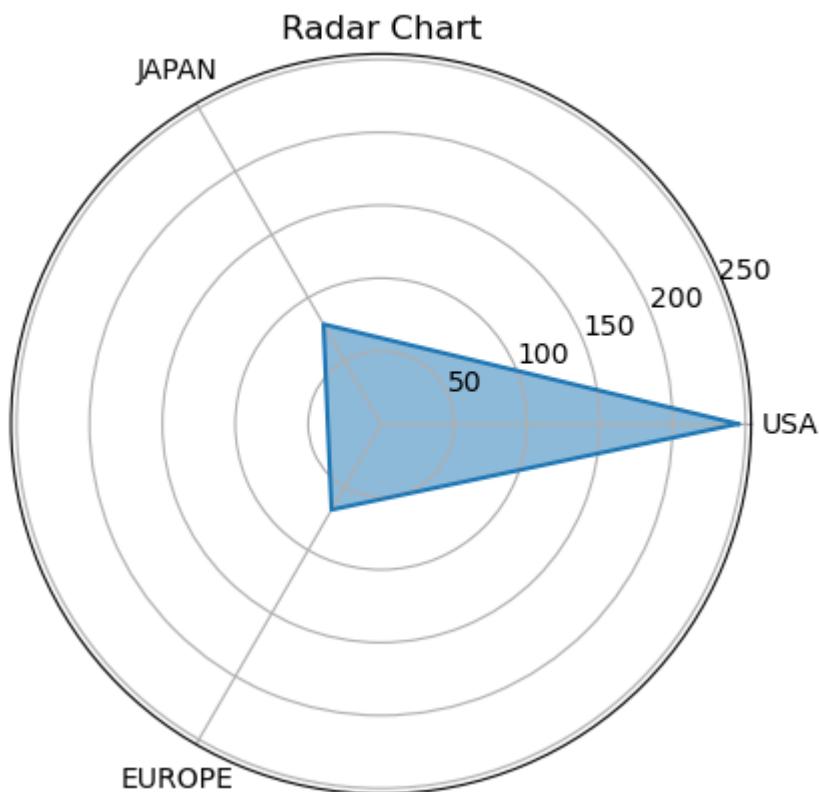
- pair plot shows plots for all attributes

Type *Markdown* and *LaTeX*:  $\alpha^2$

## Radar Chart (Spider Chart):

In [47]:

```
1 # Data for the radar chart
2 categories = ['USA', 'JAPAN', 'EUROPE']
3 values = [245, 79, 68]
4
5 # Duplicate the first value to complete the circle (closed shape)
6 values += values[:1]
7
8 # Calculate the angles for each category to plot on the polar axes
9 angles = [n / float(len(categories)) * 2 * np.pi for n in range(len(categories))]
10 angles += angles[:1]
11
12 # Create the polar plot
13 plt.polar(angles, values)
14
15 # Fill the area inside the plot with color to create a filled radar chart
16 plt.fill(angles, values, alpha=0.5)
17
18 # Set the category names as labels on the x-axis
19 plt.xticks(angles[:-1], categories)
20
21 # Add a title to the radar chart
22 plt.title('Radar Chart')
23
24 # Display the radar chart
25 plt.show()
```



## Observation

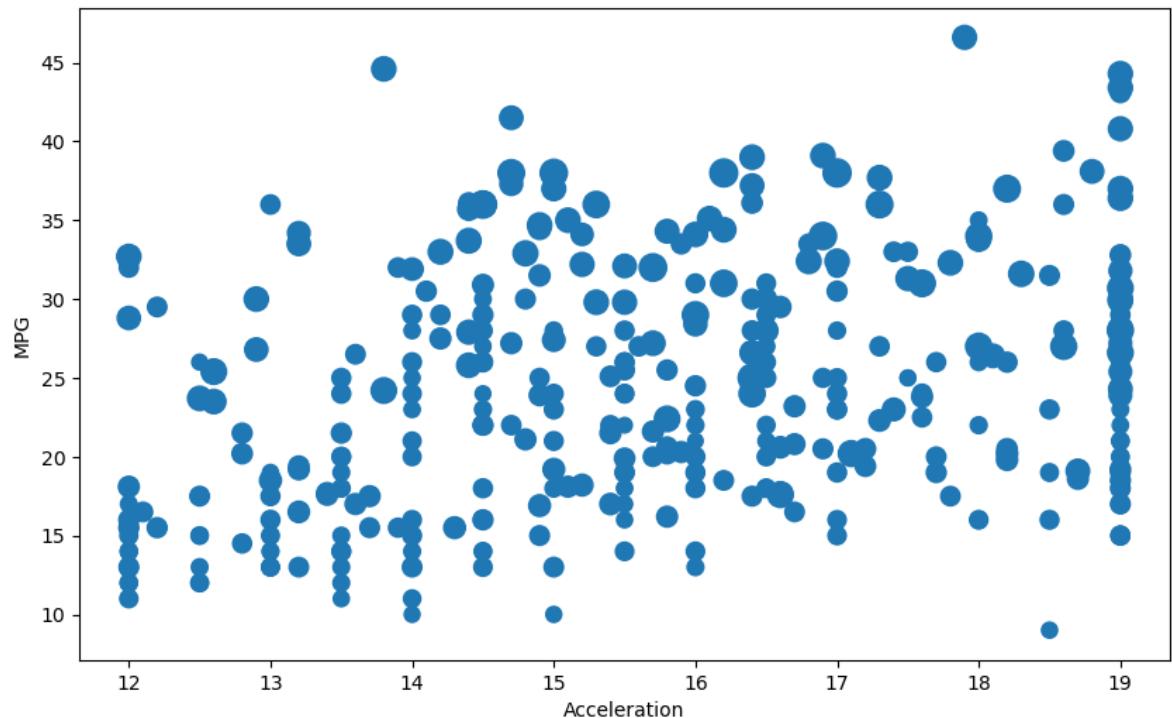
The resulting radar chart will display the values for each category, forming a closed shape with lines connecting the points for each category.

## Bubble Chart:

In [48]:

```
1 # Extract data for x-axis, y-axis, and sizes
2 x = data['acceleration']
3 y = data['mpg']
4 sizes = data['horsepower']
5
6 # Create the bubble plot with different sizes for each point
7 plt.figure(figsize=(10, 6))
8 plt.scatter(x, y, s=sizes)
9
10 # Set labels for x-axis and y-axis
11 plt.xlabel('Acceleration')
12 plt.ylabel('MPG')
13
```

Out[48]: Text(0, 0.5, 'MPG')



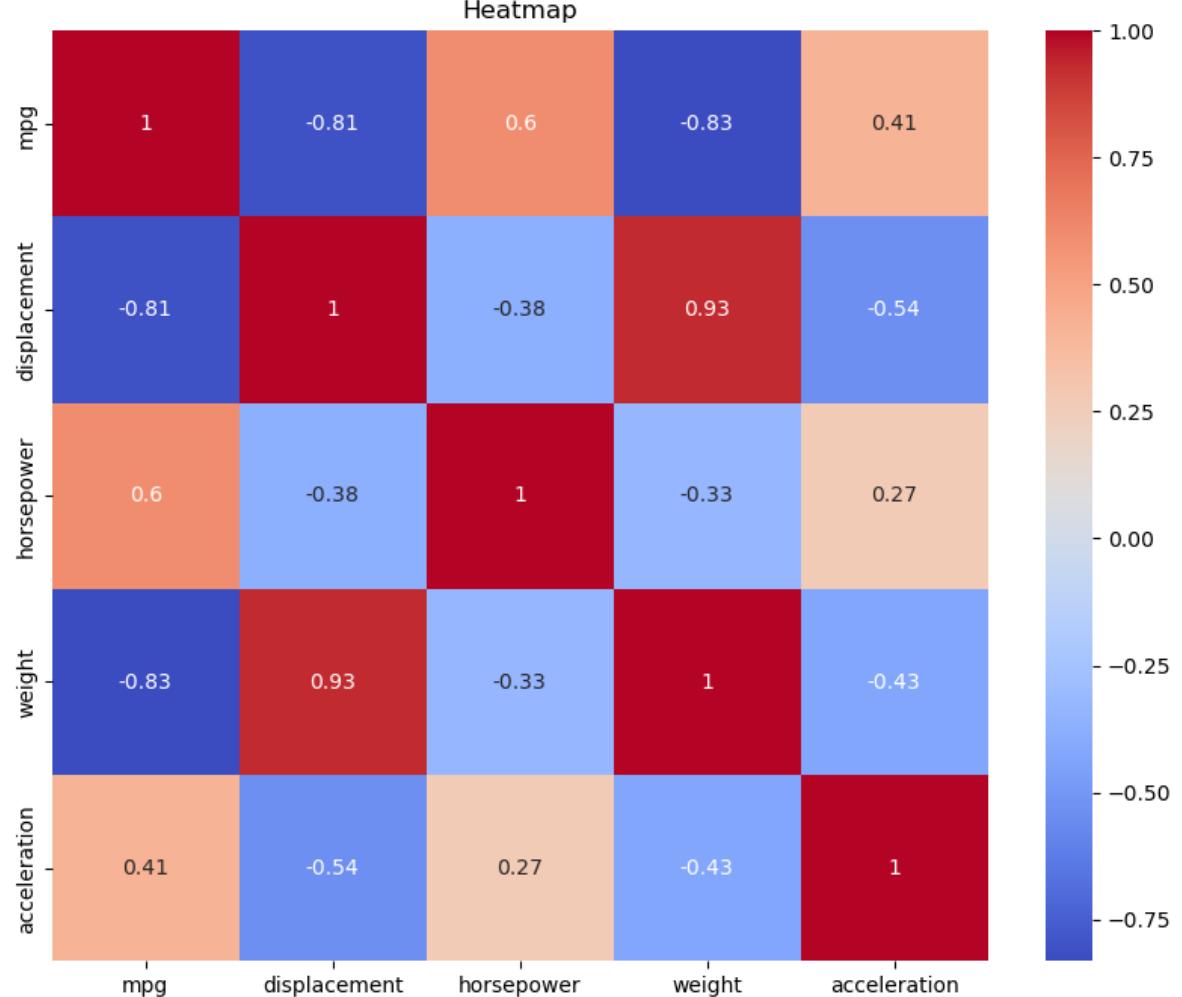
## Observation

The bubble plot is a variant of a scatter plot where the size of each data point (bubble) is determined by the 'horsepower' values from the dataset.

## Heatmap

In [49]:

```
1 # Generate the correlation matrix and create a heatmap and red represents
2 plt.figure(figsize=(10, 8))
3 sns.heatmap(data.corr(), annot=True, cmap='coolwarm')
4
5 # Set the title for the heatmap
6 plt.title('Heatmap')
7
8 # Show the heatmap
9 plt.show()
```



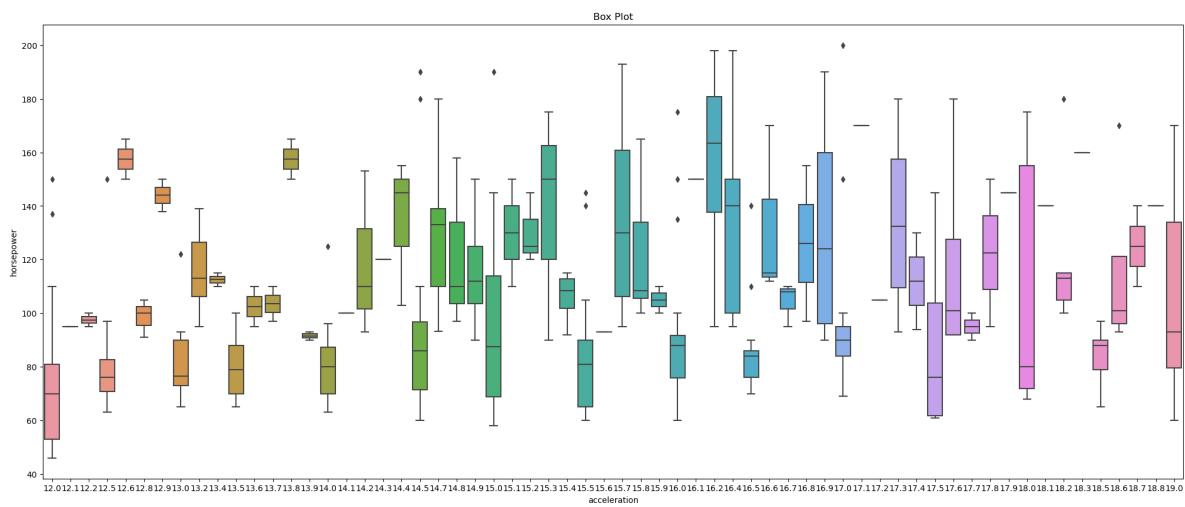
## Observation

The heatmap is a graphical representation of the correlation matrix of the dataset. It uses colors to visualize the correlation coefficients between pairs of numeric variables. 'coolwarm' colormap is used, where blue represents negative correlation, and red represents positive correlation.

## Box Plot:

In [50]:

```
1 # Create a box plot to visualize the distribution of 'horsepower' for each
2 # The 'x' axis represents the 'acceleration' values, and the 'y' axis repres
3 plt.figure(figsize=(25, 10))
4 sns.boxplot(x='acceleration', y='horsepower', data=data)
5
6 # Set the title for the box plot
7 plt.title('Box Plot')
8
9 # Show the box plot
10 plt.show()
```



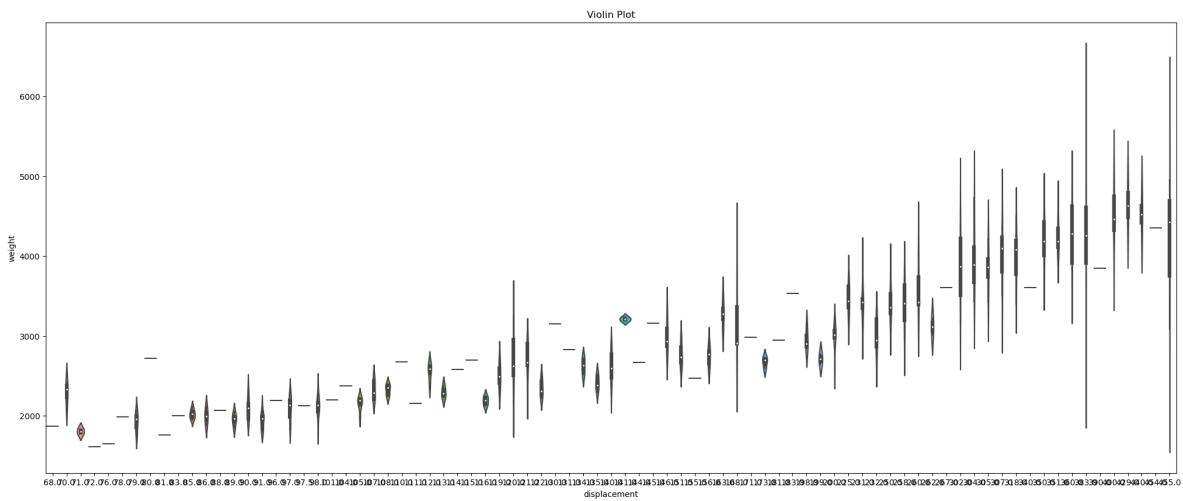
## Observation

The box plot provides a visual representation of the distribution of 'horsepower' for each 'acceleration' category. The box plot allows us to quickly assess the central tendency, variability, and presence of outliers in the 'horsepower' distribution for each 'acceleration' category.

## Violin Plot:

In [51]:

```
1 # Create a violin plot to visualize the distribution of 'weight' for each
2 # The 'x' axis represents the 'displacement' values, and the 'y' axis repri
3 plt.figure(figsize=(25, 10))
4 sns.violinplot(x='displacement', y='weight', data=data)
5
6 # Set the title for the violin plot
7 plt.title('Violin Plot')
8
9 # Show the violin plot
10 plt.show()
```



## Observation

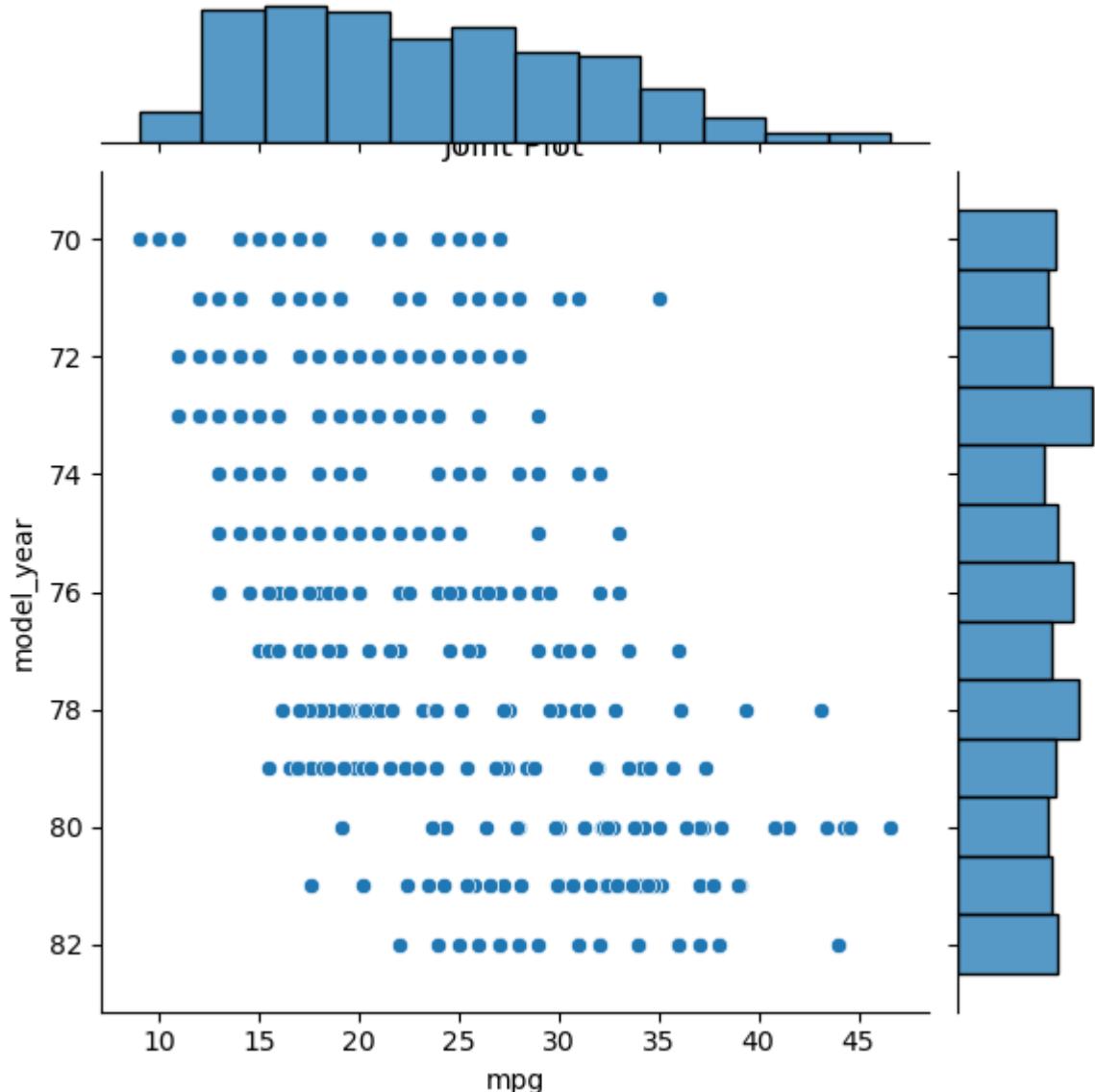
*It provides a more informative representation than a traditional box plot as it includes the kernel density plot, allowing us to observe the shape of the distribution and potential multimodality.*

## Joint Plot:

In [52]:

```
1 # Create a joint plot to visualize the relationship between 'mpg' and 'model_year'
2 plt.figure(figsize=(10, 8))
3 sns.jointplot(x='mpg', y='model_year', data=data, kind='scatter')
4
5 # Set the title for the joint plot
6 plt.title('Joint Plot')
7
8 # Show the joint plot
9 plt.show()
```

<Figure size 1000x800 with 0 Axes>



Observation

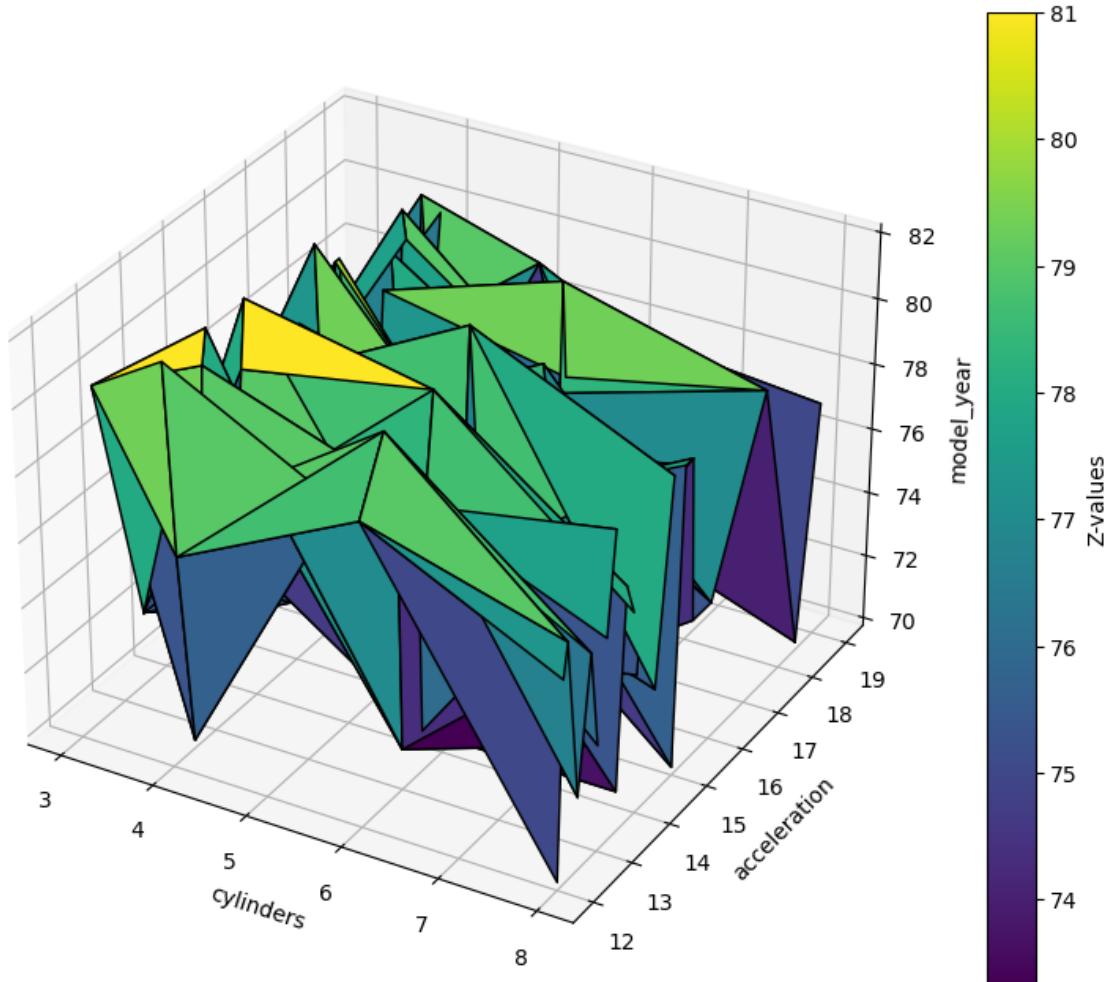
The 'y' axis represents the 'mpg' values, and the 'x' axis represents the 'model year' values

### **3D Surface Plot:**

In [53]:

```
1 # 3D surfact plot
2
3 X = data['cylinders']
4 Y = data['acceleration']
5 Z = data['model_year']
6
7 # Create the 3D surface plot
8 fig = plt.figure(figsize=(10, 8))
9 ax = fig.add_subplot(111, projection='3d')
10
11 # Plot the surface
12 ax.plot_trisurf(X, Y, Z, cmap='viridis', edgecolor='k')
13
14 # Set labels for X, Y, and Z axes
15 ax.set_xlabel('cylinders')
16 ax.set_ylabel('acceleration')
17 ax.set_zlabel('model_year')
18
19 # Add a color bar
20 fig.colorbar(ax.collections[0], label='Z-values')
21
22 plt.title('3D Surface Plot from Dataset')
23 plt.show()
24
```

3D Surface Plot from Dataset



### Observation

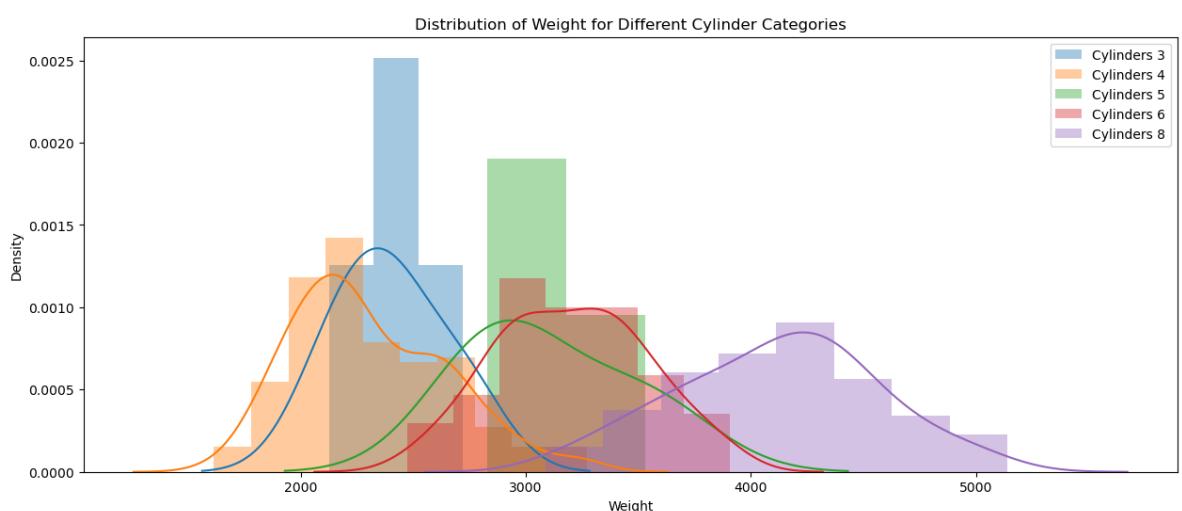
*It allows us to explore the relationship between these variables in a three-dimensional space*

Type *Markdown* and *LaTeX*:  $\alpha^2$

## **Distbution plot according to vehicles cylinder**

In [54]:

```
1 # Create a distribution plot (distplot) for 'weight' based on the 'cylinders' value
2 plt.figure(figsize=(15, 6))
3
4 # Distplot for cylinders with value 3
5 sns.distplot(data[data['cylinders'] == 3]['weight'], label='Cylinders 3')
6
7 # Distplot for cylinders with value 4
8 sns.distplot(data[data['cylinders'] == 4]['weight'], label='Cylinders 4')
9
10 # Distplot for cylinders with value 5
11 sns.distplot(data[data['cylinders'] == 5]['weight'], label='Cylinders 5')
12
13 # Distplot for cylinders with value 6
14 sns.distplot(data[data['cylinders'] == 6]['weight'], label='Cylinders 6')
15
16 # Distplot for cylinders with value 8
17 sns.distplot(data[data['cylinders'] == 8]['weight'], label='Cylinders 8')
18
19 # Add a Legend to identify each 'cylinders' category
20 plt.legend()
21
22 # Set Labels for X and Y axes
23 plt.xlabel('Weight')
24 plt.ylabel('Density')
25
26 # Set the title for the distribution plot
27 plt.title('Distribution of Weight for Different Cylinder Categories')
28
29 # Show the distribution plot
30 plt.show()
```

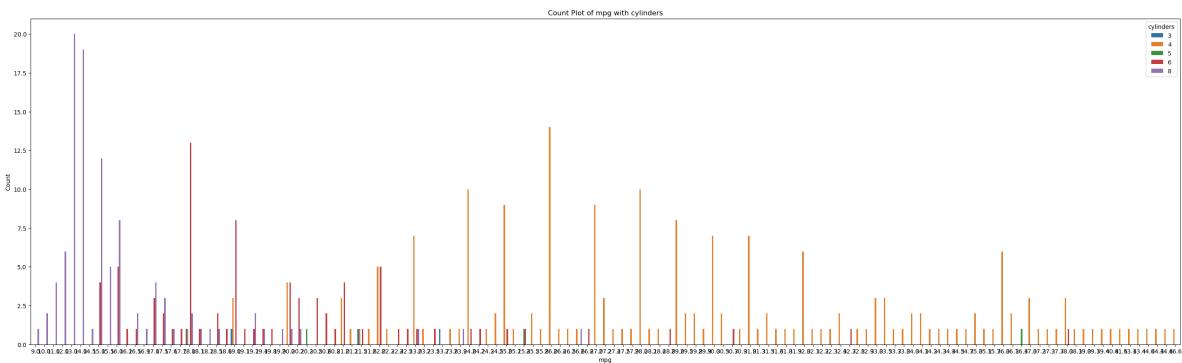


## **Observation**

The distribution plot allows us to understand how values are distributed within each category. It

In [55]:

```
1 # Create a count plot of 'mpg' with 'cylinders'
2 # The 'x' axis represents the 'mpg' values, and the count of occurrences is
3 # The 'hue' parameter is set to 'cylinders' to differentiate the count of
4 plt.figure(figsize=(35, 10))
5 sns.countplot(x='mpg', data=data, hue='cylinders')
6
7 # Set the title for the count plot
8 plt.title('Count Plot of mpg with cylinders')
9
10 # Set labels for X and Y axes
11 plt.xlabel('mpg')
12 plt.ylabel('Count')
13
14 # Show the count plot
15 plt.show()
16
17 # Crosstab to calculate count for each combination of 'mpg' and 'cylinders'
18 crosstab_result = pd.crosstab(data['cylinders'], data['mpg'])
19
20 # Display the count crosstab
21 print("Count Crosstab:")
22 print(crosstab_result)
```



Count Crosstab:

mpg	9.0	10.0	11.0	12.0	13.0	14.0	14.5	15.0	15.5	16.0	...	\
cylinders											...	
3	0	0	0	0	0	0	0	0	0	0	0	...
4	0	0	0	0	0	0	0	0	0	0	0	...
5	0	0	0	0	0	0	0	0	0	0	0	...
6	0	0	0	0	0	0	0	4	0	5	...	
8	1	2	4	6	20	19	1	12	5	8	...	

mpg	39.1	39.4	40.8	41.5	43.1	43.4	44.0	44.3	44.6	46.6	...	\
cylinders											...	
3	0	0	0	0	0	0	0	0	0	0	0	0
4	1	1	1	1	1	1	1	1	1	1	1	1
5	0	0	0	0	0	0	0	0	0	0	0	0
6	0	0	0	0	0	0	0	0	0	0	0	0
8	0	0	0	0	0	0	0	0	0	0	0	0

[5 rows x 127 columns]

## Observation

*The count bars are color-coded based on the 'cylinders' categories. Each bar represents the count of cars with a specific number of cylinders for the corresponding 'mpg' value. The count crosstab provides a tabular representation of the counts for each combination of 'mpg' and 'cylinders.'*