

```
In [9]: import pandas as pd  
import numpy as np
```

```
In [10]: bike = pd.read_csv("D:/dataset/Bike/biketrain.csv", parse_dates=["datetime"])  
titanic = pd.read_csv("D:/dataset/titanic_data/train_modified.csv")  
print(bike.shape, titanic.shape)  
  
(10886, 12) (891, 12)
```

```
In [11]: # !pip install brewer2mpl
```

Collecting brewer2mpl

Downloading <https://files.pythonhosted.org/packages/84/57/00c45a199719e617db0875181134fcb3aeef701deae346547ac722eaa5e/brewer2mpl-1.4.1-py2.py3-none-any.whl> (<https://files.pythonhosted.org/packages/84/57/00c45a199719e617db0875181134fcb3aeef701deae346547ac722eaa5e/brewer2mpl-1.4.1-py2.py3-none-any.whl>)

Installing collected packages: brewer2mpl

Successfully installed brewer2mpl-1.4.1

```
In [12]: import numpy as np
import pandas as pd
import matplotlib as mpl
import matplotlib.pyplot as plt
import seaborn as sns
import warnings; warnings.filterwarnings(action='once')

large = 22; med = 16; small = 12
params = {'axes.titlesize': large,
          'legend.fontsize': med,
          'figure.figsize': (16, 10),
          'axes.labelsize': med,
          'axes.titlesize': med,
          'xtick.labelsize': med,
          'ytick.labelsize': med,
          'figure.titlesize': large}
plt.rcParams.update(params)
plt.style.use('seaborn-whitegrid')
sns.set_style("white")
%matplotlib inline

# Version
print(mpl.__version__) #> 3.0.0
print(sns.__version__) #> 0.9.0
```

2.1.2

0.8.1

Correlation

In [4]: bike.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10886 entries, 0 to 10885
Data columns (total 12 columns):
datetime      10886 non-null datetime64[ns]
season        10886 non-null int64
holiday       10886 non-null int64
workingday    10886 non-null int64
weather       10886 non-null int64
temp          10886 non-null float64
atemp        10886 non-null float64
humidity      10886 non-null int64
windspeed     10886 non-null float64
casual        10886 non-null int64
registered    10886 non-null int64
count         10886 non-null int64
dtypes: datetime64[ns](1), float64(3), int64(8)
memory usage: 1020.6 KB
```

In [5]: titanic.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 12 columns):
PassengerId   891 non-null int64
Survived      891 non-null int64
Pclass        891 non-null int64
Name          891 non-null object
Sex           891 non-null object
Age           891 non-null float64
SibSp         891 non-null int64
Parch         891 non-null int64
Ticket        891 non-null object
Fare          891 non-null float64
Cabin         204 non-null object
Embarked      891 non-null object
dtypes: float64(2), int64(5), object(5)
memory usage: 83.6+ KB
```

01. Scatter Plot

- 2개 또는 그 이상의 변수의 관계를 확인

```
In [26]: # Import Dataset
df = pd.read_csv("https://github.com/selva86/datasets/raw/master/mtcars.csv")
print(df.shape)
print(df.head())
print(df.columns)
print(df.info())
```

```
(32, 14)
```

	mpg	cyl	disp	hp	drat	wt	qsec	vs	am	gear	carb	fast	W
0	4.582576	6	160.0	110	3.90	2.620	16.46	0	1	4	4	1	
1	4.582576	6	160.0	110	3.90	2.875	17.02	0	1	4	4	1	
2	4.774935	4	108.0	93	3.85	2.320	18.61	1	1	4	1	1	
3	4.626013	6	258.0	110	3.08	3.215	19.44	1	0	3	1	1	
4	4.324350	8	360.0	175	3.15	3.440	17.02	0	0	3	2	1	

	cars	carname
0	Mazda RX4	Mazda RX4
1	Mazda RX4 Wag	Mazda RX4 Wag
2	Datsun 710	Datsun 710
3	Hornet 4 Drive	Hornet 4 Drive
4	Hornet Sportabout	Hornet Sportabout

```
Index(['mpg', 'cyl', 'disp', 'hp', 'drat', 'wt', 'qsec', 'vs', 'am', 'gear',
      'carb', 'fast', 'cars', 'carname'],
      dtype='object')
```

```
<class 'pandas.core.frame.DataFrame'>
```

```
RangeIndex: 32 entries, 0 to 31
```

```
Data columns (total 14 columns):
```

```
mpg      32 non-null float64
```

```
cyl      32 non-null int64
```

```
disp     32 non-null float64
```

```
hp       32 non-null int64
```

```
drat     32 non-null float64
```

```
wt       32 non-null float64
```

```
qsec     32 non-null float64
```

```
vs       32 non-null int64
```

```
am       32 non-null int64
```

```
gear     32 non-null int64
```

```
carb     32 non-null int64
```

```
fast     32 non-null int64
```

```
cars     32 non-null object
```

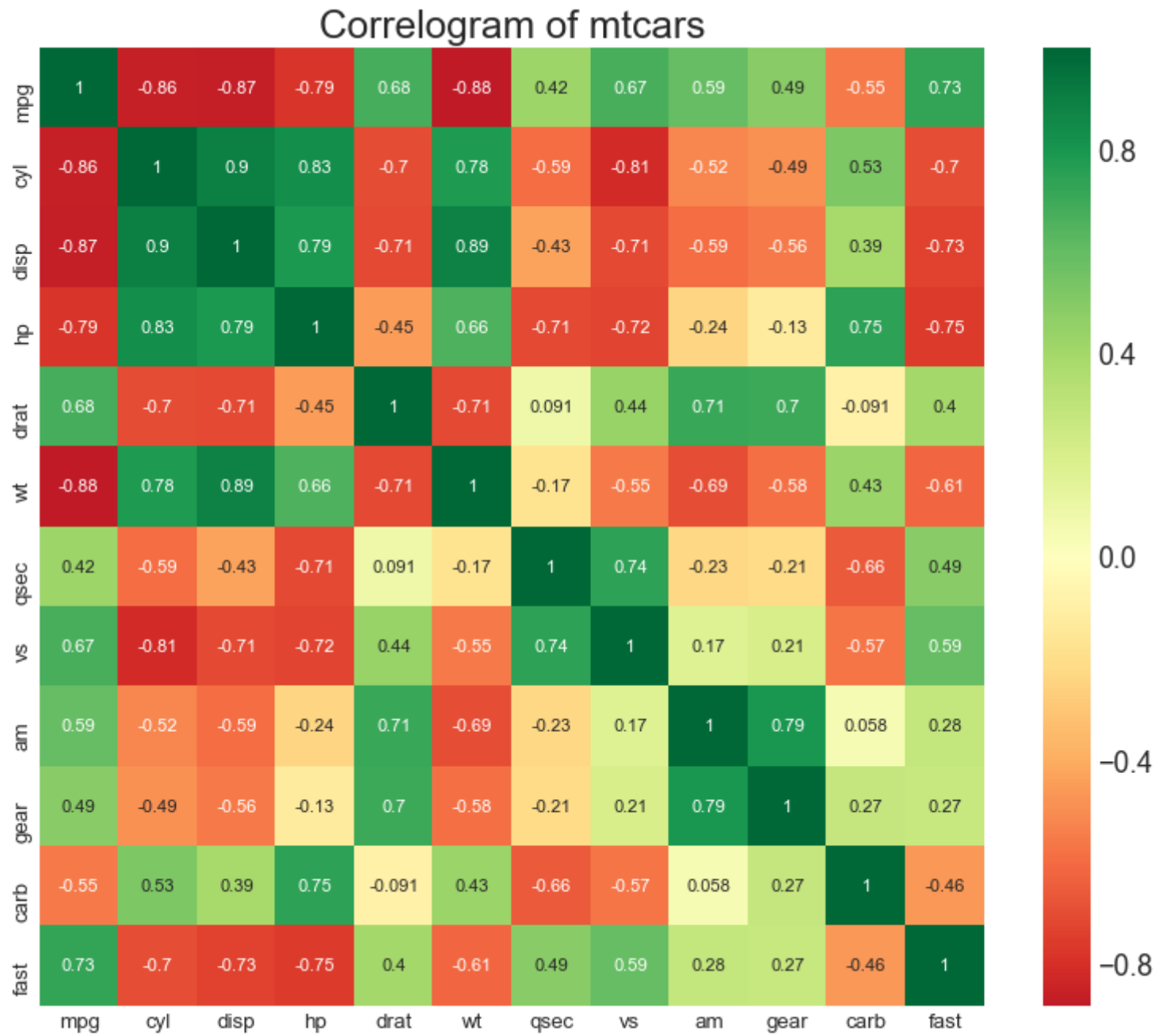
```
carname  32 non-null object
```

```
dtypes: float64(5), int64(7), object(2)  
memory usage: 3.6+ KB  
None
```

In []:

```
In [24]: # Plot
plt.figure(figsize=(12,10), dpi= 80)
sns.heatmap(df.corr(),
            xticklabels=df.corr().columns,
            yticklabels=df.corr().columns, cmap='RdYlGn', center=0, annot=True)

# Decorations
plt.title('Correlogram of mtcars', fontsize=22)
plt.xticks(fontsize=12)
plt.yticks(fontsize=12)
plt.show()
```



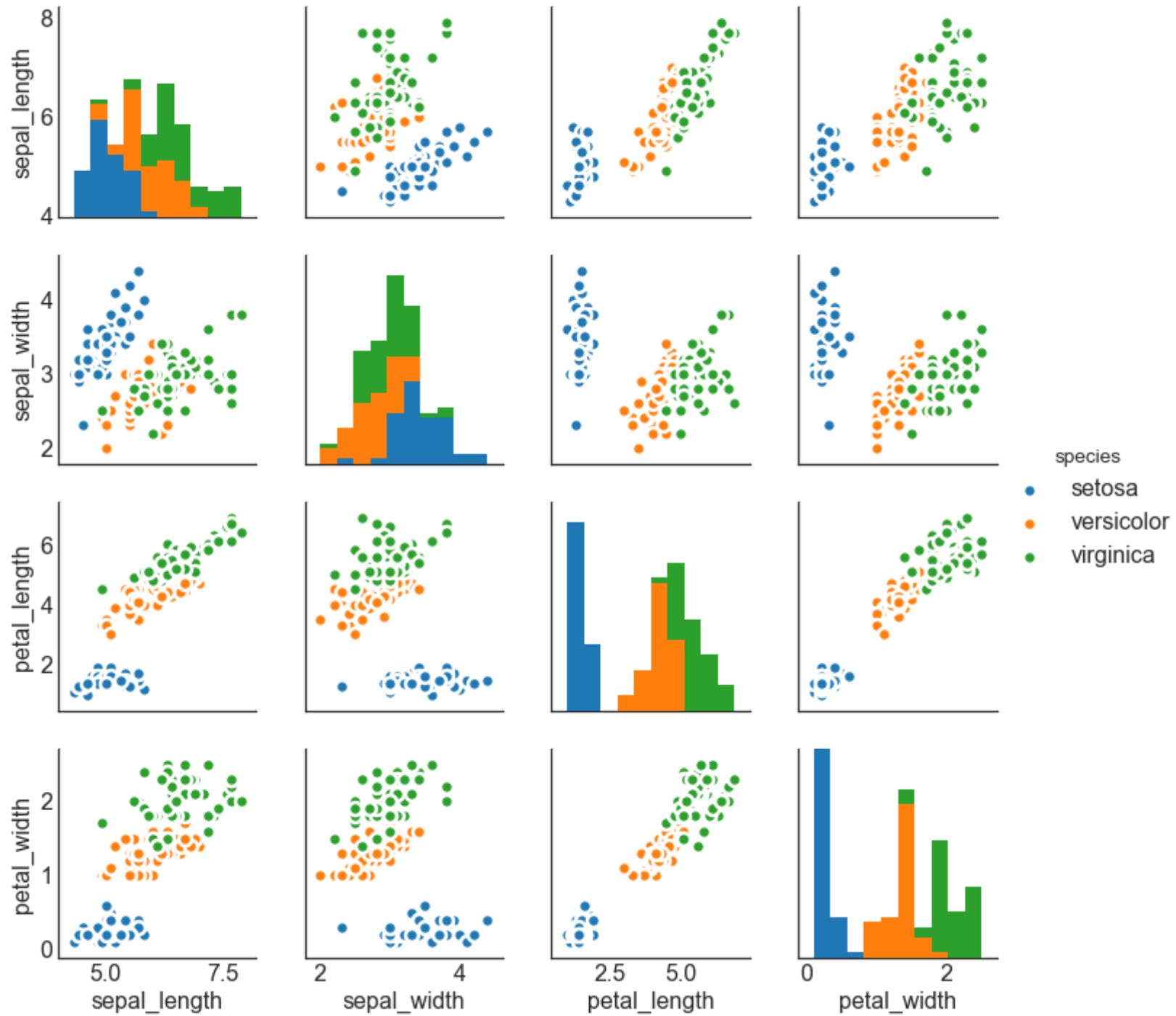
Pairwise Plot

- Pairwise Plot is a favorite in exploratory analysis to understand the relationship between all possible pairs of numeric variables. It is a must have tool for bivariate analysis
- 수치형 변수의 가능한 모든 쌍 사이의 관계를 이해하기 위해 좋다.

```
In [27]: # Load Dataset
df = sns.load_dataset('iris')

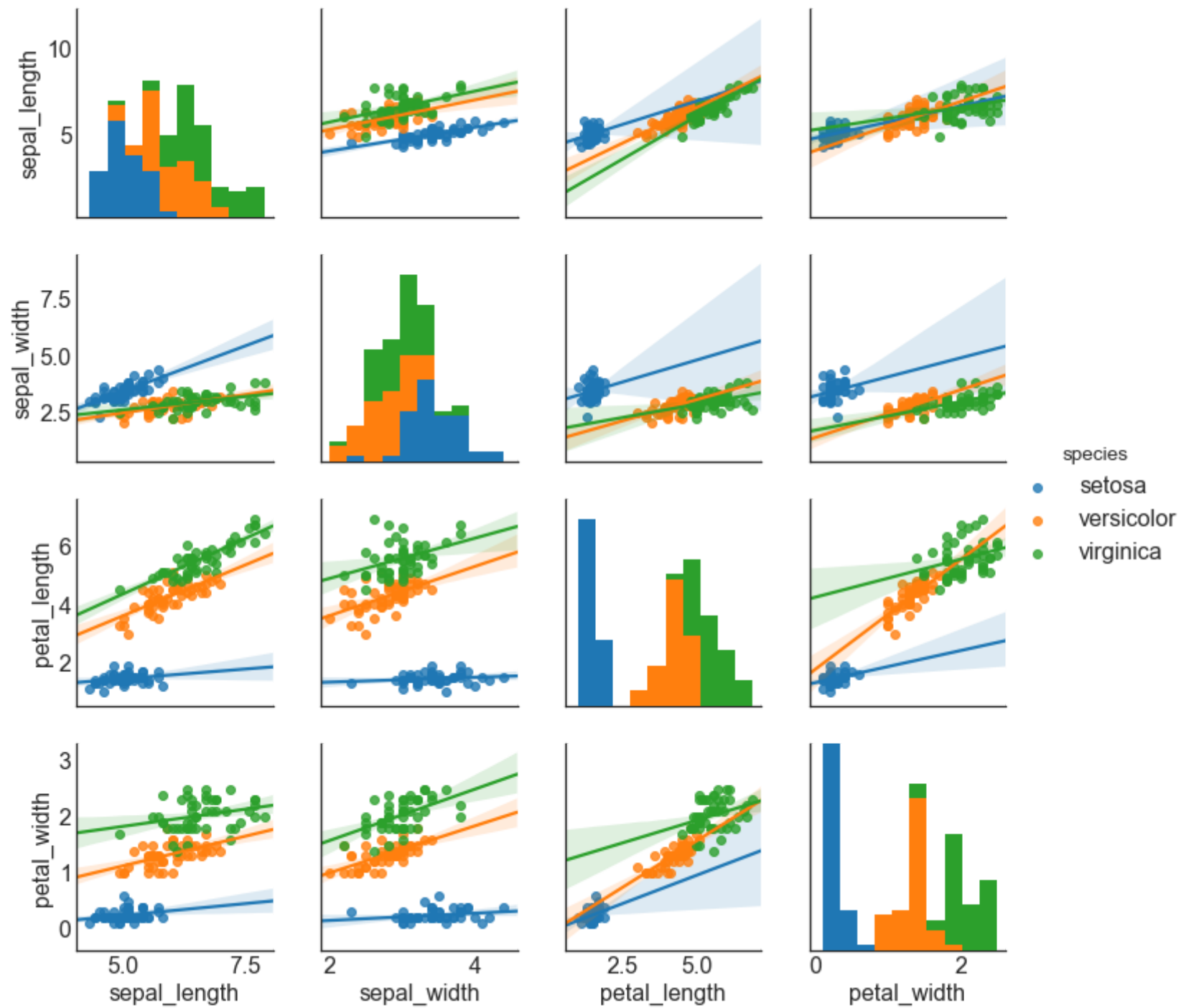
# Plot
plt.figure(figsize=(10,8), dpi= 80)
sns.pairplot(df, kind="scatter", hue="species", plot_kws=dict(s=80, edgecolor="white", linewidth=2.5))
plt.show()
```

<matplotlib.figure.Figure at 0x1e919f8ccf8>




```
In [28]: # Load Dataset  
df = sns.load_dataset('iris')  
  
# Plot  
plt.figure(figsize=(10,8), dpi= 80)  
sns.pairplot(df, kind="reg", hue="species")  
plt.show()
```

<matplotlib.figure.Figure at 0x1e919f41a58>



15. Ordered Bar Chart

Fuel economy data from 1999 and 2008 for 38 popular models of car
A data frame with 234 rows and 11 variables

This dataset contains a subset of the fuel economy data that the EPA makes available on <http://fuelconomy.gov>. It contains only models which had a new release every year between 1999 and 2008 – this was used as a proxy for the popularity of the car.

이 데이터 세트에는 EPA가 <http://fuelconomy.gov>에서 제공하는 연비 데이터의 일부를 포함하고 있습니다. 1999 년에서 2008 년 사이에 매년 새 버전이 출시 된 모델 만 포함되어 있으며 이는 자동차의 인기를 위한 프록시로 사용되었습니다.

manufacturer model : model name
displ : engine displacement, in litres
year : year of manufacture
cyl : number of cylinders
trans : type of transmission
drv : f = front-wheel drive, r = rear wheel drive, 4 = 4wd
cty : city miles per gallon (갤런당 도시 마일)
hwy : highway miles per gallon
fl : fuel type
class : "type" of car

```
In [29]: # Import Dataset
df_raw = pd.read_csv("https://github.com/selva86/datasets/raw/master/mpg_ggplot2.csv")
print(df_raw.shape)
print(df_raw.head())
print(df_raw.columns)
print(df_raw.info())
```

```
(234, 11)
  manufacturer model  displ  year  cyl    trans drv  cty  hwy fl   class
0         audi   a4    1.8  1999    4  auto(l5)  f   18   29  p  compact
1         audi   a4    1.8  1999    4 manual(m5)  f   21   29  p  compact
2         audi   a4    2.0  2008    4 manual(m6)  f   20   31  p  compact
3         audi   a4    2.0  2008    4  auto(av)   f   21   30  p  compact
4         audi   a4    2.8  1999    6  auto(l5)   f   16   26  p  compact
Index(['manufacturer', 'model', 'displ', 'year', 'cyl', 'trans', 'drv', 'cty',
      'hwy', 'fl', 'class'],
      dtype='object')
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 234 entries, 0 to 233
Data columns (total 11 columns):
manufacturer    234 non-null object
model           234 non-null object
displ           234 non-null float64
year           234 non-null int64
cyl            234 non-null int64
trans          234 non-null object
drv            234 non-null object
cty            234 non-null int64
hwy            234 non-null int64
fl             234 non-null object
class          234 non-null object
dtypes: float64(1), int64(4), object(6)
memory usage: 20.2+ KB
None
```



```
In [38]: df = df_raw[['cty', 'manufacturer']].groupby('manufacturer').apply(lambda x: x.mean())  
df
```

Out[38]:

cty	
manufacturer	
audi	17.611111
chevrolet	15.000000
dodge	13.135135
ford	14.000000
honda	24.444444
hyundai	18.642857
jeep	13.500000
land rover	11.500000
lincoln	11.333333
mercury	13.250000
nissan	18.076923
pontiac	17.000000
subaru	19.285714
toyota	18.529412
volkswagen	20.925926

```
In [49]: df.index
```

Out[49]: Index(['audi', 'chevrolet', 'dodge', 'ford', 'honda', 'hyundai', 'jeep',
 'land rover', 'lincoln', 'mercury', 'nissan', 'pontiac', 'subaru',
 'toyota', 'volkswagen'],
 dtype='object', name='manufacturer')

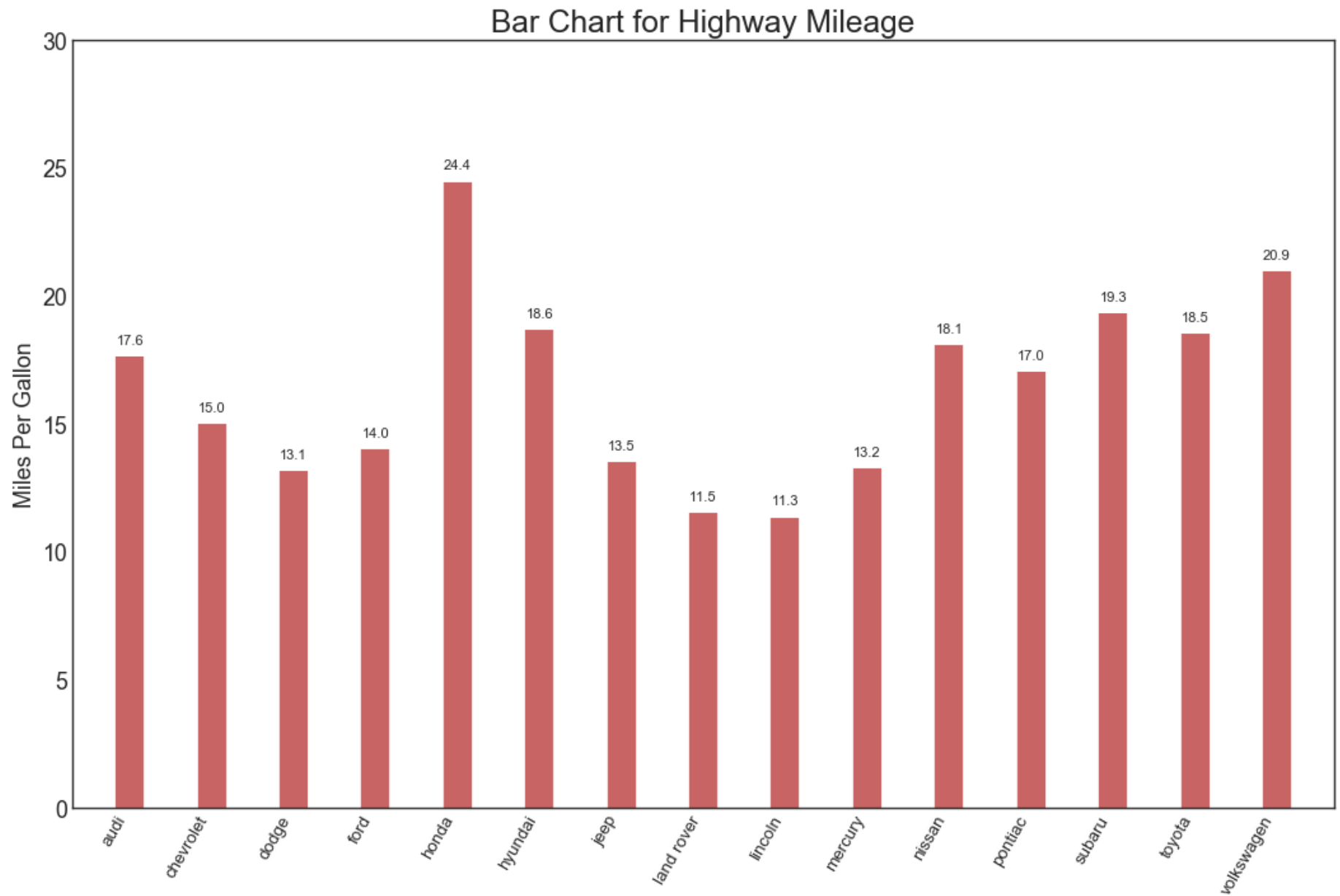
```
In [52]: # Draw plot
import matplotlib.patches as patches

fig, ax = plt.subplots(figsize=(16,10), facecolor='white', dpi= 80)
ax.vlines(x=df.index, ymin=0, ymax=df.cty, color='firebrick', alpha=0.7, linewidth=20)

# Annotate Text
for i, cty in enumerate(df.cty):
    ax.text(i, cty+0.5, round(cty, 1), horizontalalignment='center')

# Title, Label, Ticks and Ylim
ax.set_title('Bar Chart for Highway Mileage', fontdict={'size':22})
ax.set_ylabel('Miles Per Gallon', ylim=(0, 30))
plt.xticks(df.index, rotation=60, horizontalalignment='right', fontsize=12)

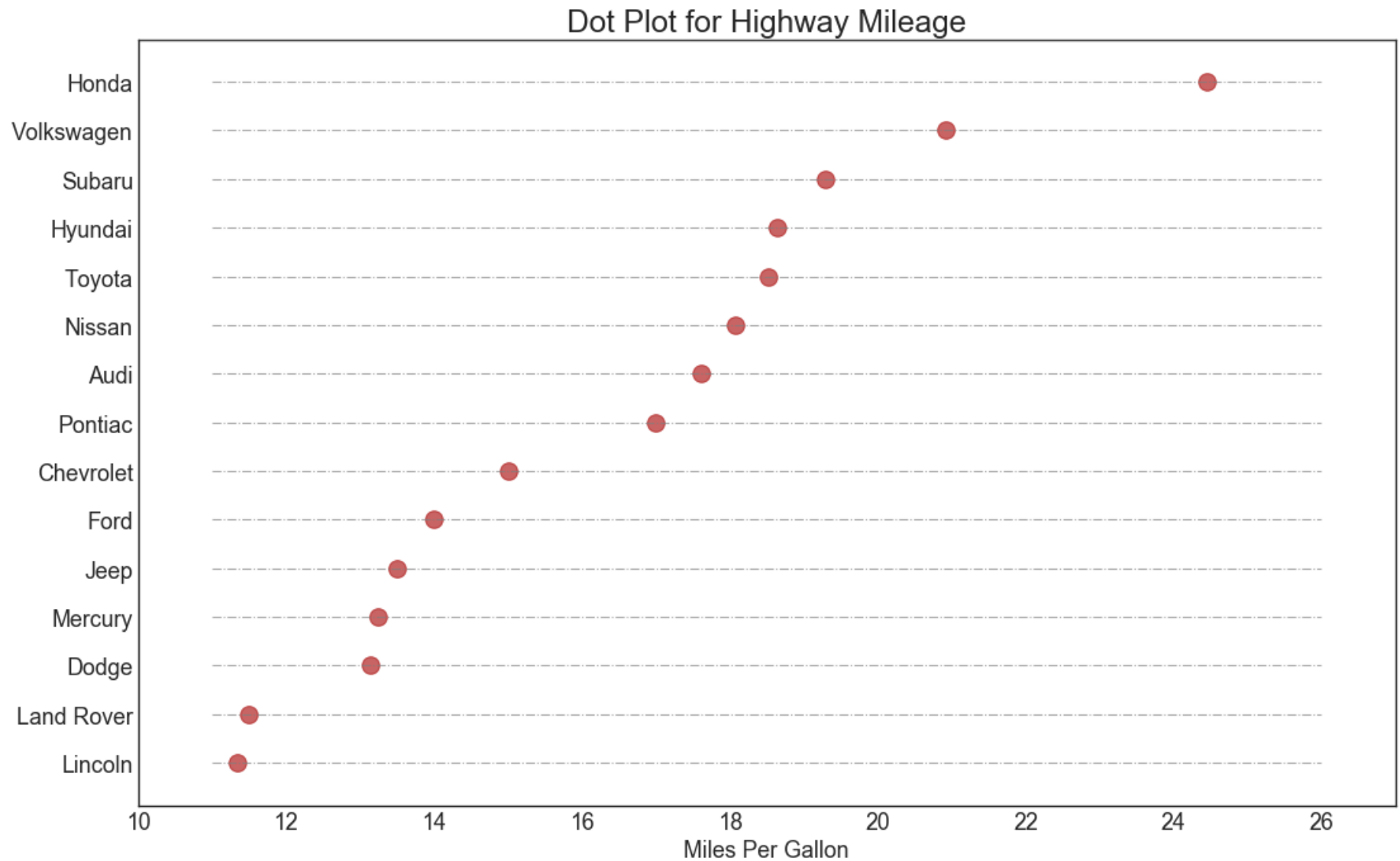
plt.show()
```



```
In [59]: df = df_raw[['cty', 'manufacturer']].groupby('manufacturer').apply(lambda x: x.mean())
df.sort_values('cty', inplace=True)
df.reset_index(inplace=True)

# Draw plot
fig, ax = plt.subplots(figsize=(16,10), dpi= 80)
ax.hlines(y=df.index, xmin=11, xmax=26, color='gray', alpha=0.7, linewidth=1, linestyle='dashdot')
ax.scatter(y=df.index, x=df.cty, s=150, color='firebrick', alpha=0.7)

# Title, Label, Ticks and Ylim
ax.set_title('Dot Plot for Highway Mileage', fontdict={'size':22})
ax.set_xlabel('Miles Per Gallon')
ax.set_yticks(df.index)
ax.set_yticklabels(df.manufacturer.str.title(),
                    fontdict={'horizontalalignment': 'right'})
ax.set_xlim(10, 27)
plt.show()
```



Density Plot

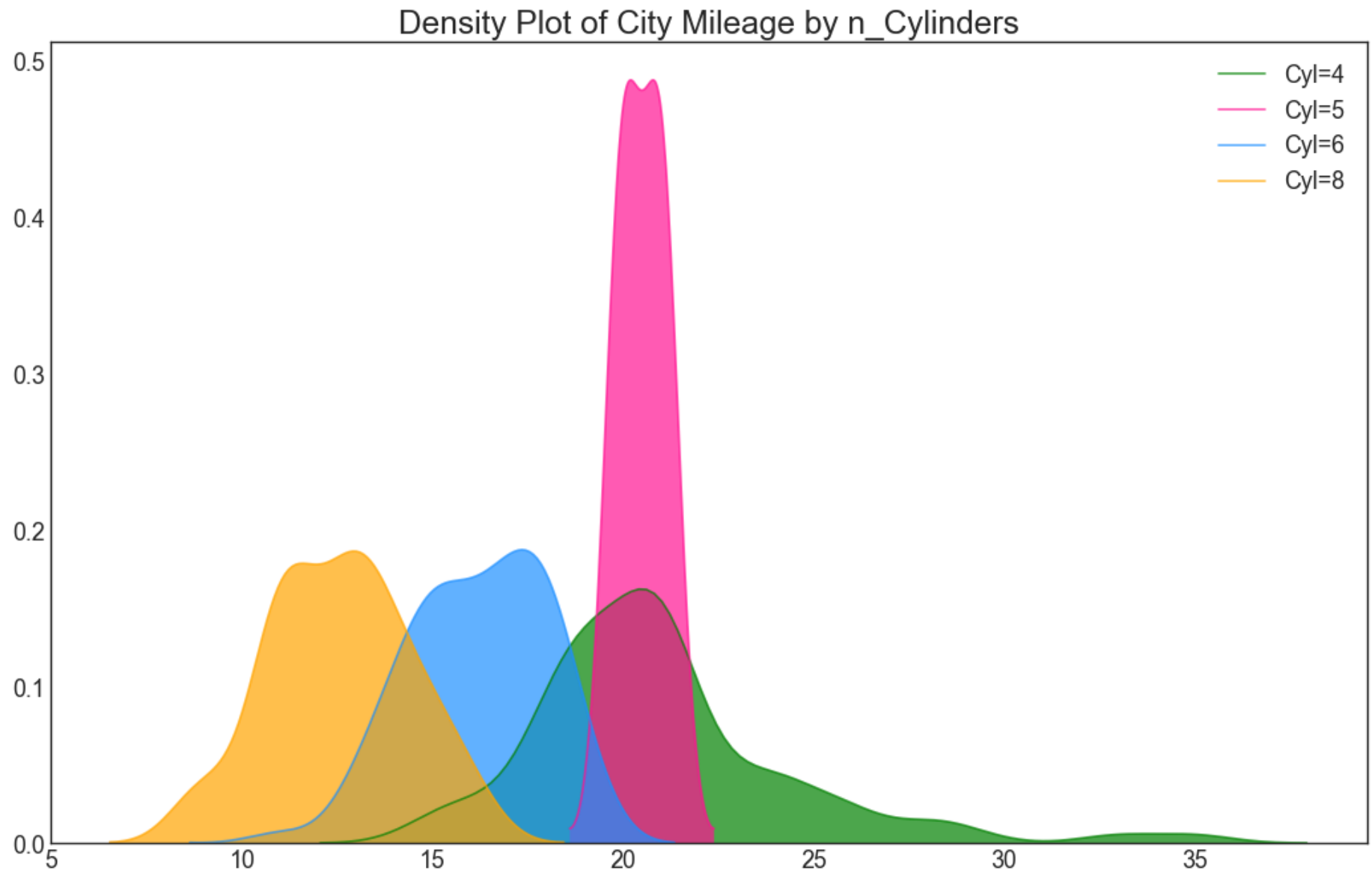
- 밀도 플롯은 일반적으로 사용되는 도구로 연속 변수의 분포를 시각화합니다. '응답' 변수로 그룹화하면 X와 Y의 관계를 검사 할 수 있습니다. 도시 마일리지 분포가 실린더 수에 따라 어떻게 달라지는지를 표현하기위한 목적이 아래의 경우입니다.

```
In [60]: # Import Data
df = pd.read_csv("https://github.com/selva86/datasets/raw/master/mpg_ggplot2.csv")

# Draw Plot
plt.figure(figsize=(16,10), dpi= 80)
sns.kdeplot(df.loc[df['cyl'] == 4, "cty"], shade=True, color="g", label="Cyl=4", alpha=.7)
sns.kdeplot(df.loc[df['cyl'] == 5, "cty"], shade=True, color="deeppink", label="Cyl=5", alpha=.7)
sns.kdeplot(df.loc[df['cyl'] == 6, "cty"], shade=True, color="dodgerblue", label="Cyl=6", alpha=.7)
sns.kdeplot(df.loc[df['cyl'] == 8, "cty"], shade=True, color="orange", label="Cyl=8", alpha=.7)

# Decoration
plt.title('Density Plot of City Mileage by n_Cylinders', fontsize=22)
plt.legend()
plt.show()
```

C:\Users\WWITHJSW\Anaconda3\lib\site-packages\statsmodels\nonparametric\kde.py:475: DeprecationWarning: object of type <class 'numpy.float64'> cannot be safely interpreted as an integer.
grid,delta = np.linspace(a,b,gridsize,retstep=True)



```
In [65]: !pip install squarify
```

Requirement already satisfied: squarify in c:\Users\Wwithjs\Anaconda3\lib\site-packages (0.3.0)

C:\Users\WWITHJS\Anaconda3\lib\site-packages\IPython\utils\process_win32.py:131: ResourceWarning: unclosed file <_io.BufferedWriter name=6>

return process_handler(cmd, _system_body)

C:\Users\WWITHJS\Anaconda3\lib\site-packages\IPython\utils\process_win32.py:131: ResourceWarning: unclosed file <_io.BufferedReader name=7>

return process_handler(cmd, _system_body)

C:\Users\WWITHJS\Anaconda3\lib\site-packages\IPython\utils\process_win32.py:131: ResourceWarning: unclosed file <_io.BufferedReader name=8>

return process_handler(cmd, _system_body)

33. Treemap

- 트리 맵은 파이차트와 유사하다. 각 그룹의 기여도를 과도하지 않고 잘 보여준다.


```
In [70]: import squarify
# Import Data
df_raw = pd.read_csv("https://github.com/selva86/datasets/raw/master/mpg_ggplot2.csv")

# Prepare Data
df = df_raw.groupby('class').size().reset_index(name='counts')
labels = df.apply(lambda x: str(x[0]) + "Wn (" + str(x[1]) + ")", axis=1)
sizes = df['counts'].values.tolist()
colors = [plt.cm.Spectral(i/float(len(labels))) for i in range(len(labels))]

# Draw Plot
plt.figure(figsize=(16,12), dpi= 90)
squarify.plot(sizes=sizes, label=labels, color=colors, alpha=.8)

# Decorate
plt.title('Treemap of Vechile Class')
plt.axis('off')
plt.show()
```

Treemap of Vechile Class



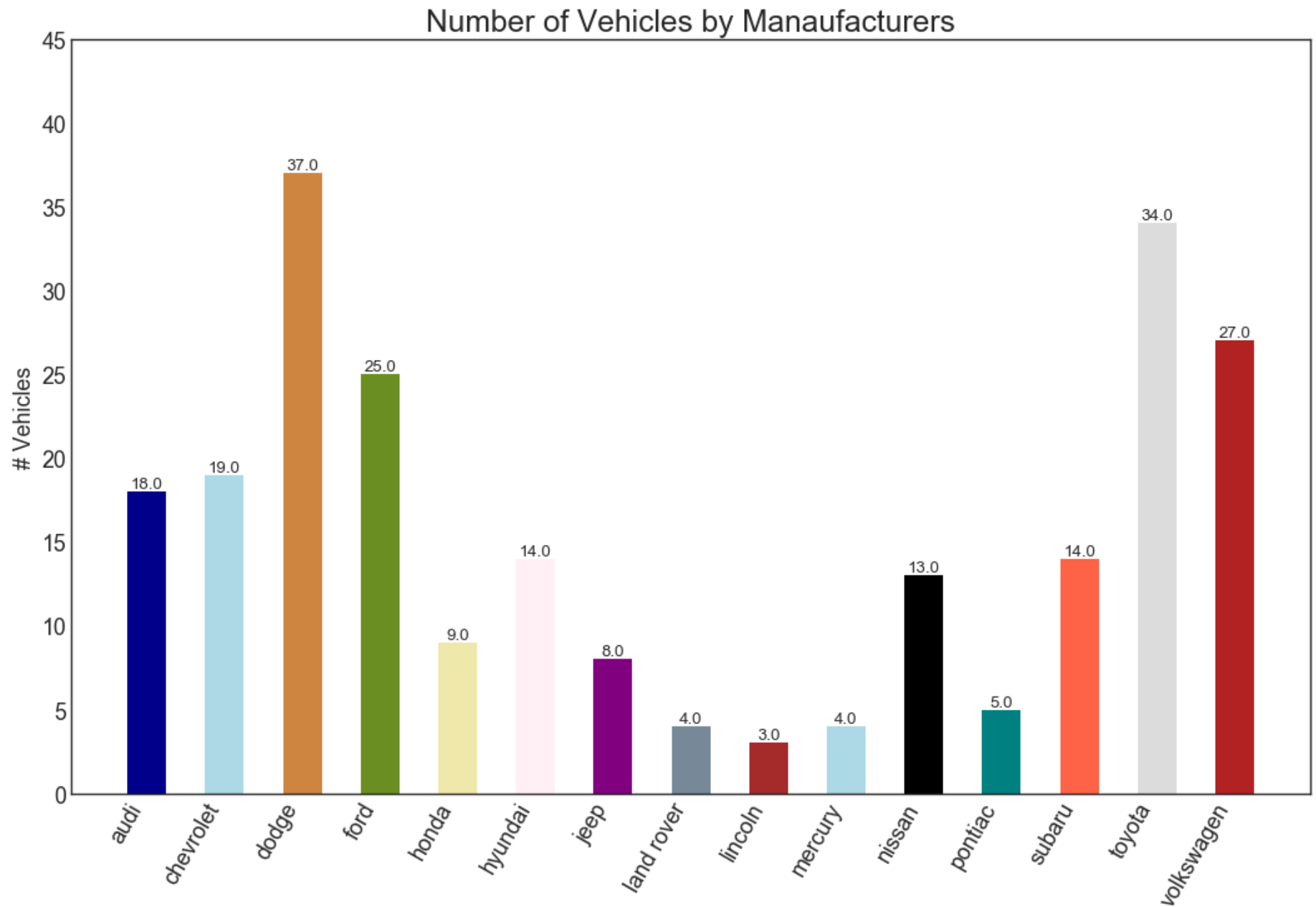
```
In [71]: import random

# Import Data
df_raw = pd.read_csv("https://github.com/selva86/datasets/raw/master/mpg_ggplot2.csv")

# Prepare Data
df = df_raw.groupby('manufacturer').size().reset_index(name='counts')
n = df['manufacturer'].unique().__len__()+1
all_colors = list(plt.cm.colors.cnames.keys())
random.seed(100)
c = random.choices(all_colors, k=n)

# Plot Bars
plt.figure(figsize=(16,10), dpi= 80)
plt.bar(df['manufacturer'], df['counts'], color=c, width=.5)
for i, val in enumerate(df['counts'].values):
    plt.text(i, val, float(val), horizontalalignment='center', verticalalignment='bottom', fontdict={'fontweight':500, 'size':12})

# Decoration
plt.gca().set_xticklabels(df['manufacturer'], rotation=60, horizontalalignment= 'right')
plt.title("Number of Vehicles by Manufacturers", fontsize=22)
plt.ylabel('# Vehicles')
plt.ylim(0, 45)
plt.show()
```

**REF**

<https://www.machinelearningplus.com/plots/top-50-matplotlib-visualizations-the-master-plots-python/>
(<https://www.machinelearningplus.com/plots/top-50-matplotlib-visualizations-the-master-plots-python/>)

In []: