

In [1]:

```
# This Python 3 environment comes with many helpful analytics libraries installed
# It is defined by the kaggle/python docker image: https://github.com/kaggle/docker-python
# For example, here's several helpful packages to load in

import numpy as np # linear algebra
import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)

# Input data files are available in the "../input/" directory.
# For example, running this (by clicking run or pressing Shift+Enter) will list all files under the input directory

import os
for dirname, _, filenames in os.walk('../kaggle/input'):
    for filename in filenames:
        print(os.path.join(dirname, filename))

# Any results you write to the current directory are saved as output.
```

/kaggle/input/car-mpg/mpg_raw.csv

This is a 3 part series in which i will walk through a dataset analysing it and then at the end do predictive modelling. I recommend to follow the parts in sequential order but you can jump to any part by clicking on the heading links in blue.

Part 1, Exploratory Data Analysis(EDA):
This part consists of summary statistics of data but the major focus will be on EDA where we extract meaning/information from data using plots and report important insights about data. This part is more about **data analysis and business intelligence(BI)**.

Part 2,Statistical Analysis:
In this part we will do many statistical hypothesis testing, apply estimation statistics and interpret the results we get. We will also validate this with the findings from part one. We will apply both parametric and non-parametric tests. We will report all the important insights we get in this part. This part is all about **data science** requires statistical background.

Part 3, Predictive Modelling:
In this part we will predict some response using predictors. This part is all about **machine learning**.

If you like these notebooks then please upvote and also share with others.

Data Description

The data we are using for EDA is the [auto mpg](#) dataset taken from UCI repository.

Information regarding data

Title: Auto-Mpg Data

Number of Attributes: 9 including the class attribute

Attribute Information:

1. mpg: continuous

2. cylinders: multi-valued discrete

3. displacement: continuous

4. horsepower: continuous

5. weight: continuous

6. acceleration: continuous

7. model year: multi-valued discrete

8. origin: multi-valued discrete

9. car name: string (unique for each instance)

All the attributes are self-explanatory.

This data is not complex and is good for analysis as it has a nice blend of both categorical and numerical attributes.

[data source](#)

This is **part i.e.** EDA. We won't stretch this part long and do following things in sequential manner.

1. **Preprocess the data**, this includes dealing with missing values, duplicate data if any and then align the data.

2. **EDA on categorical attributes**, this includes analysing their distributions and relations with other cat. attributes.

3. **EDA on numerical attributes**, this includes analysing their distributions and relations with other num. attributes.

4. Then we will analyse the **relation b/w num. & cat. attributes**.

I make use of **seaborn** heavily throughout the notebook, so it is also a good goto notebook for those who are looking for EDA using seaborn.

In [2]:

```
# First import all necessary libraries
import itertools
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
```

In [3]:

```
# set seaborn's default settings
sns.set()
```

We will first import the data into a pandas dataframe and inspect it's properties.

In [4]:

```
df = pd.read_csv("../input/car-mpg/mpg_raw.csv")
df.head()
```

Out[4]:

	mpg	cylinders	displacement	horsepower	weight	acceleration	model_year	origin	name
0	18.0	8	307.0	130.0	3504	12.0	70	usa	chevrolet chevelle malibu
1	15.0	8	350.0	165.0	3693	11.5	70	usa	buick skylark 320
2	18.0	8	318.0	150.0	3436	11.0	70	usa	plymouth satellite
3	16.0	8	204.0	150.0	3433	12.0	70	usa	amc rebel sst
4	17.0	8	302.0	140.0	3449	10.5	70	usa	ford torino

In [5]:

```
# so now the data is in rectangular form with 398 entries each having 9 distinct properties
df.shape
```

Out[5]: (398, 9)

In [6]:

```
# Lets list all the columns
columns = list(df.columns)
columns
```

Out[6]:

['mpg', 'cylinders', 'displacement', 'horsepower', 'weight', 'acceleration', 'model_year', 'origin', 'name']

In [7]:

```
# we now describe the properties of this dataframe like column datatype etc.
df.info()
```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 398 entries, 0 to 397
Data columns (total 9 columns):
mpg 398 non-null float64
cylinders 398 non-null int64
displacement 398 non-null float64
horsepower 392 non-null float64
weight 398 non-null int64
acceleration 398 non-null float64
model_year 398 non-null object
origin 398 non-null object
dtypes: float64(4), int64(3), object(2)
memory usage: 28.1+ KB

We now make two distinct list for categorical and numerical column names as the analysis differ for both the types. For that we introspect the datatypes of each column and if it is of type 'object' then it's categorical and else numerical.

We will use these two lists heavily throughout the analysis.

In [8]:

```
cats = list(df.select_dtypes(include=['object']).columns)
nums = list(df.select_dtypes(exclude=['object']).columns)
print(f'categorical variables: {cats}')
print(f'numerical variables: {nums}')
```

categorical variables: ['origin', 'name']
numerical variables: ['mpg', 'cylinders', 'displacement', 'horsepower', 'weight', 'acceleration', 'model_year']

In [9]:

```
# Let's inspect how many unique values are there in each column.
df.nunique(axis=0)
```

Out[9]:

mpg 129
cylinders 5
displacement 82
horsepower 83
weight 351
acceleration 95
model_year 13
origin 3
name 385
dtype: int64

As there are very few unique values for cylinders and model_year, so it's safe to make them categorical instead of numeric.

In [10]:

```
# cylinders and model_year also seems to be categorical so lets update the lists
cats.extend(['cylinders', 'model_year'])
nums.remove('cylinders')
nums.remove('model_year')
```

print(f'categorical variables: {cats}')
print(f'numerical variables: {nums}')

categorical variables: ['origin', 'name', 'cylinders', 'model_year']
numerical variables: ['mpg', 'displacement', 'horsepower', 'weight', 'acceleration']

Now inspect for nans in data.

In [11]:

```
# check for 'nans' in each column
df.isna().sum()
```

Out[11]:

mpg 0
cylinders 0
displacement 0
horsepower 6
weight 0
acceleration 0
model_year 0
origin 0
name 0
dtype: int64

In [12]:

```
# Let's print these 6 'nan' containing rows
df[df.isnull().any(axis=1)]
```

Out[12]:

	mpg	cylinders	displacement	horsepower	weight	acceleration	model_year	origin	name
32	25.0	4	98.0	NaN	2046	19.0	71	usa	ford pinto
126	21.0	6	200.0	NaN	1875	17.0	74	usa	ford maverick
330	40.9	4	85.0	NaN	1835	17.3	80	europa	renault lecar deluxe
336	23.6	4	140.0	NaN	2905	14.3	81	usa	ford mustang cobra
354	34.5	4	100.0	NaN	2320	15.8	80	europa	renault 18i
374	23.0	4	151.0	NaN	3035	20.5	82	usa	amc concord dl

In [13]:

```
# nan rows proportion in data
6 / len(df)
```

Out[13]: 0.01587537688442211

So horsepower consists of total of 6 nan rows comprising of around 1.5% of data. As this fraction is very low so it's safe to drop these nan rows for now.

Note: If the nan-row proportion is large enough then we won't drop it but instead impute missing values.

In [14]:

```
# for now remove all nan rows as they are just 1.5%
df = df.dropna(inplace=True)
df.reset_index(inplace=True, axis=1)
df.shape
```

Out[14]: (392, 9)

In [15]:

```
# find total duplicate entries and drop them if any
print(f'total duplicate rows: {df.duplicated().sum()}')
```

drop duplicate rows if any
df = df[df.duplicated()]
df.shape

total duplicate rows: 0

Out[15]: (392, 9)

In [16]:

```
# before we move ahead it's a good practice to group all variables together having same type.
df = pd.concat([df[cats], df[nums]], axis=1)
df.head()
```

Out[16]:

	origin	name	cylinders	model_year	mpg	displacement	horsepower	weight	acceleration
0	usa	chevrolet chevelle malibu	8	70	18.0	307.0	130.0	3504	12.0
1	usa	buick skylark 320	8	70	15.0	350.0	165.0	3693	11.5
2	usa	plymouth satellite	8	70	18.0	318.0	150.0	3436	11.0
3	usa	amc rebel sst	8	70	16.0	204.0	150.0	3433	12.0
4	usa	ford torino	8	70	17.0	302.0	140.0	3449	10.5

In [17]: num_rows, num_cols = df.shape

In [18]:

```
# save this cleaned df to csv
df.to_csv('mpg_cleaned.csv', index=False)
```

Now we're all good to go for some in-depth analysis

Analysis on Categorical Attributes

Our analysis includes both descriptive stats and EDA.

In [19]:

```
# Let's import the cleaned version of mpg although no need here because we already updated df
df = pd.read_csv('mpg_cleaned.csv')
```

In [20]:

```
print(f'categorical variables: {cats}')
```

categorical variables: ['origin', 'name', 'cylinders', 'model_year']

We will first slice out the categorical columns from original dataframe and then do analysis on it keeping the original data untouched, and at the end incorporate needed changes in our original dataframe.

In [21]:

```
df_cat = df.loc[:, 'origin':'model_year']
df_cat.head()
```

Out[21]:

	origin	name	cylinders	model_year
0	usa	chevrolet chevelle malibu	8	70
1	usa	buick skylark 320	8	70
2	usa	plymouth satellite	8	70
3	usa	amc rebel sst	8	70
4	usa	ford torino	8	70

As origin and name consists of text data so it needs some preprocessing. We will remove all extra spaces from each string. otherwise the same string with different spacings will be treated as different categories which should not be the case.

In [22]:

```
# remove extra spaces if any
for col in ['origin', 'name']:
    df_cat[col] = df_cat[col].apply(lambda x: ' '.join(x.split()))
```

We will create an artificial categorical attribute named 'mpg_level' which categorizes mpg into low, medium and high. This is done for two reasons, first it will help a lot in EDA i.e. we can bifurcate plots on the basis of mpg and secondly this is easy to understand as compared to numbers in mpg.

We are dividing mpg into three regions as,
(min,17]> low
(17,29]> medium
(29,max]> high

Also the choice of the range is analytical and can be anything till it seems to be reasonable.

Note: This is feature-engineering and mostly done in predictive modeling but it makes sense to introduce it here.

In [23]:

```
df_cat['mpg_level'] = df['mpg'].apply(lambda x: 'low' if x<17 else 'high' if x>=29 else 'medium')
cats.append('mpg_level')
```

print(f'categorical variables: {cats}')

categorical variables: ['origin', 'name', 'cylinders', 'model_year', 'mpg_level']

In [24]:

```
# Let's look at the unique categories in 'origin', 'cylinders' & 'model_year'
# we are leaving name because it is almost unique for each entry (nothing interesting)
print(f'categories in origin: {pd.unique(df_cat["origin"])}')
print(f'categories in cylinders: {pd.unique(df_cat["cylinders"])}')
print(f'categories in model_year: {pd.unique(df_cat["model_year"])}')
```

categories in origin: ['usa', 'japan', 'europe']
categories in cylinders: [8 4 6 3 5]
categories in model_year: [70 71 72 73 74 75 76 77 78 79 80 81 82]

In [25]:

```
# Although descriptive stats for categorical variables are not much informative but still it's worth
# Also pandas describe function is only for numeric data and in df_cat 'cylinders' & 'model_year' are
# the only numeric type.
df_cat.describe()
```

Out[25]:

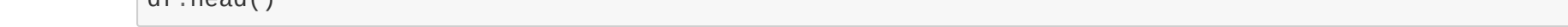
	cylinders	model_year
count	392.000000	392.000000
mean	5.471939	75.979592
std	1.705783	3.686373
min	3.000000	70.000000
25%	4.000000	73.000000
50%	4.000000	76.000000
75%	8.000000	79.000000
max	8.000000	82.000000

It seems that most of the values in 'cylinders' are 4 with a min of 3 and max of 8.

Analysis of Distribution

Now we analyse the distribution for each categorical feature and make some insights from the plots.

In case of categorical variables an ideal (or atleast loved) distribution is **uniform**.



In [26]:

```
fig = plt.figure(1, (14, 8))
```

for i,cat in enumerate(df_cat.drop(['name'], axis=1).columns):
ax = plt.subplot(2,2,i+1)
sns.countplot(df_cat[cat], order=df_cat[cat].value_counts().index)
ax.set_xlabel(None)
ax.set_title(f'Distribution of {cat}')
plt.tight_layout()

plt.show()



In [27]:

```
# calculate proportion of dominant classes in each category
for i,cat in enumerate(df_cat.drop(['name'], axis=1).columns):
    val_counts = df_cat[cat].value_counts()
    dominant_frac = val_counts.iloc[0] / num_rows
    print(f'val_counts.index[0] : num_rows : dominant_frac = {round(dominant_frac * 100, 2)}% of {cat}')
```

'usa' alone contributes to 62.5% of origin
'4' alone contributes to 56.7% of cylinders
'73' alone contributes to 18.2% of model_year
'medium' alone contributes to 52.3% of mpg_level

In [28]:

```
# count of different cylinders
df_cat.cylinders.value_counts()
```

Out[28]:

4 199
8 193
6 83
3 34
5 3
Name: cylinders, dtype: int64

Insights

• origin is highly imbalanced, usa alone consists of 62.5% of data whereas japan & europe are having similar proportion.

We will see this dominance in future analysis. We will try to find the reason for this in our further analysis.

• cylinders is highly imbalanced, 4 alone consists of 50.7% of data. Whereas 8 & 6 are nearly in same proportion but 3 & 5 collectively accounts for only 7 entries i.e. 1.8% of entire data. We will see this huge proportional imbalance in our future analysis.

• mpg_level is highly imbalanced, medium alone consists of 52.3% of data while low & high are in the same proportion. This dominance is due the fact of our thresholding while manufacturing this feature because the medium range is broader hence it consists of more data points. It won't be there in original 'mpg' feature as it is continuous.

• model_year is considerably balanced which is good.

Now we analyse car 'name'.

In [29]:

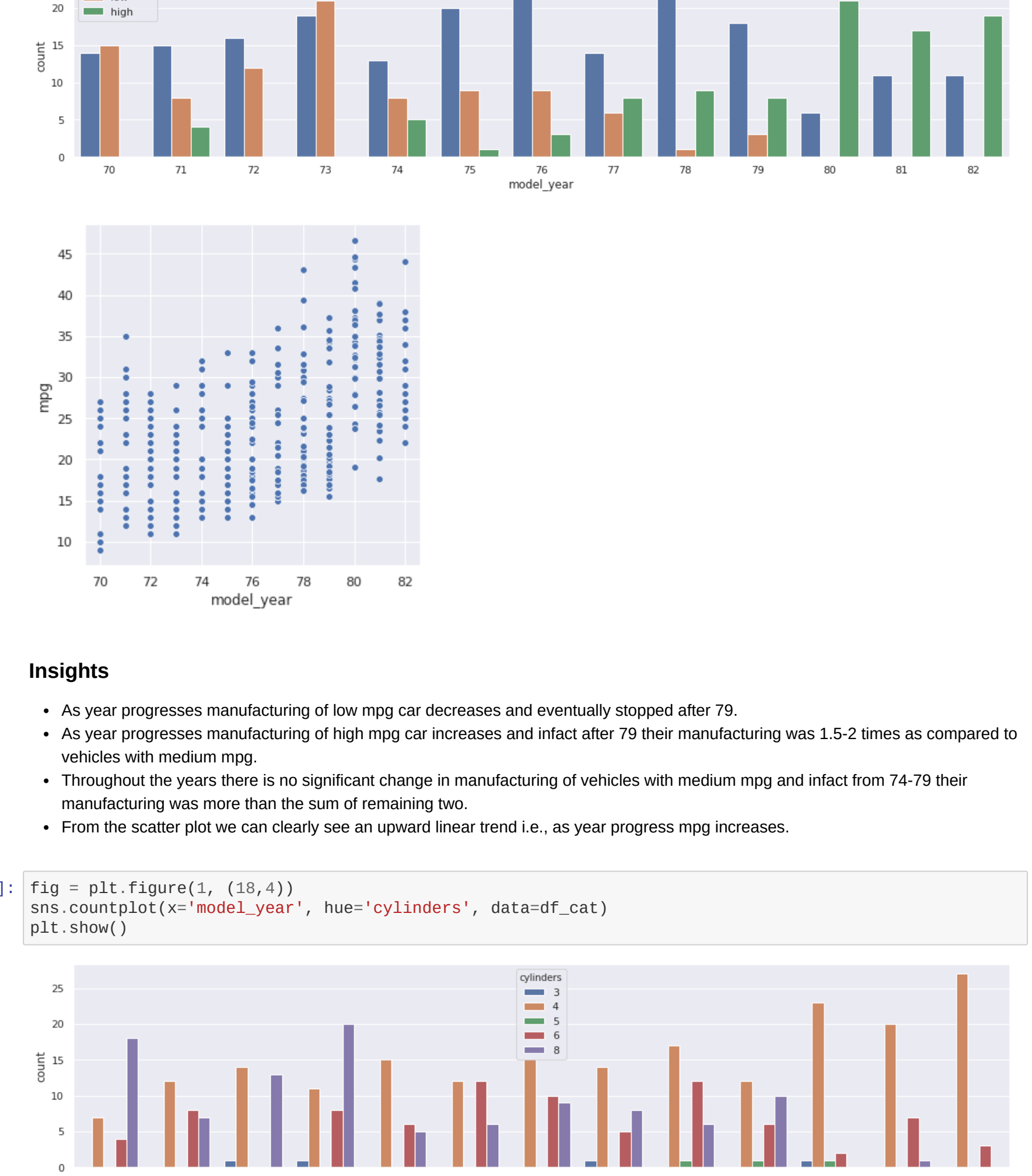
```
print(f'total unique categories in 'name': {df_cat.name.nunique()}')
print(f'unique categories in 'name':\n\n{df_cat.name.unique()}')
```

total unique categories in 'name': 381

unique categories in 'name':

['chevrolet chevelle malibu', 'buick skylark 320', 'plymouth satellite', 'amc rebel sst', 'ford torino', 'ford galaxie 500', 'chevrolet impala', 'plymouth fury ii', 'pontiac catalina', 'amc ambassador dpl', 'dodge challenger se', 'plymouth cuda 340', 'chevrolet monte carlo', 'dodge challenger se', 'toyota corona mark ii', 'plymouth duster', 'amc hornet', 'ford naverick', 'datsun p180i', 'volkswagen 113i deluxe sedan', 'peugeot 504', 'audi 100 ls', 'saab 990', 'bmw 280i', 'amc gremlin', 'ford f250', 'chevy c20', 'dodge d200', 'hi 1200d', 'chevrolet Vega 2300', 'toyota corona', 'plymouth custom', 'audi fox', 'ford torino 500', 'amc matador', 'pontiac catalina brougham', 'dodge d300', 'dodge monaco', 'ford country squire (sw)', 'pontiac safari (sw)', 'amc hornet sportabout (sw)', 'chevrolet Vega (sw)', 'pontiac firebird', 'ford mustang', 'mercury capri 2000', 'opel 1900', 'peugeot 304', 'fiat 124b', 'toyota corolla 1200', 'datsun 1200', 'volkswagen model 111', 'toyota corona mark ii (sw)', 'dodge colt (sw)', 'chevrolet malibu', 'volkswagen 113i deluxe sedan', 'peugeot 504', 'audi 100 ls', 'saab 990', 'bmw 280i', 'amc gremlin', 'ford f250', 'chevy c20', 'dodge d200', 'hi 1200d', 'chevrolet Vega 2300', 'toyota corona', 'plymouth custom', 'audi fox', 'ford torino 500', 'amc matador', 'pontiac catalina brougham', 'dodge d300', 'dodge monaco', 'ford country squire (sw)', 'pontiac safari (sw)', 'amc hornet sportabout (sw)', 'chevrolet Vega (sw)', 'pontiac firebird', 'ford mustang', 'mercury capri 2000', 'opel 1900', 'peugeot 304', 'fiat 124b', 'toyota corolla 1200', 'datsun 1200', 'volkswagen model 111', 'toyota corona mark ii (sw)', 'dodge colt (sw)', 'chevrolet malibu', 'volkswagen 113i deluxe sedan', 'peugeot 504', 'audi 100 ls', 'saab 990', 'bmw 280i', 'amc gremlin', 'ford f250', 'chevy c20', 'dodge d200', 'hi 1200d', 'chevrolet Vega 2300', 'toyota corona', 'plymouth custom', 'audi fox', 'ford torino 500', 'amc matador', 'pontiac catalina brougham', 'dodge d300', 'dodge monaco', 'ford country squire (sw)', 'pontiac safari (sw)', 'amc hornet sportabout (sw)', 'chevrolet Vega (sw)', 'pontiac firebird', 'ford mustang', 'mercury capri 2000', 'opel 1900', 'peugeot 304', 'fiat 124b', 'toyota corolla 1200', 'datsun 1200', 'volkswagen model 111', 'toyota corona mark ii (sw)', 'dodge colt (sw)', 'chevrolet malibu', 'volkswagen 113i deluxe sedan', 'peugeot 504', 'audi 100 ls', 'saab 990', 'bmw 280i', 'amc gremlin', 'ford f250', 'chevy c20', 'dodge d200', 'hi 1200d', 'chevrolet Vega 2300', 'toyota corona', 'plymouth custom', 'audi fox', 'ford torino 500', 'amc matador', 'pontiac catalina brougham', 'dodge d300', 'dodge monaco', 'ford country squire (sw)', 'pontiac safari (sw)', 'amc hornet sportabout (sw)', 'chevrolet Vega (sw)', 'pontiac firebird', 'ford mustang', 'mercury capri 2000', 'opel 1900', 'peugeot 304', 'fiat 124b', 'toyota corolla 1200', 'datsun 1200', 'volkswagen model 111', 'toyota corona mark ii (sw)', 'dodge colt (sw)', 'chevrolet malibu', 'volkswagen 113i deluxe sedan', 'peugeot 504', 'audi 100 ls', 'saab 990', 'bmw 280i', 'amc gremlin', 'ford f250', 'chevy c20', 'dodge d200', 'hi 1200d', 'chevrolet Vega 2300', 'toyota corona', 'plymouth custom', 'audi fox', 'ford torino 500', 'amc matador', 'pontiac catalina brougham', 'dodge d300', 'dodge monaco', 'ford country squire (sw)', 'pontiac safari (sw)', 'amc hornet sportabout (sw)', 'chevrolet Vega (sw)', 'pontiac firebird', 'ford mustang', 'mercury capri 2000', 'opel 1900', 'peugeot 304', 'fiat 124b', 'toyota corolla 1200', 'datsun 1200', 'volkswagen model 111', 'toyota corona mark ii (sw)', 'dodge colt (sw)', 'chevrolet malibu', 'volkswagen 113i deluxe sedan', 'peugeot 504', 'audi 100 ls', 'saab 990', 'bmw 280i', 'amc gremlin', 'ford f250', 'chevy c20', 'dodge d200', 'hi 1200d', 'chevrolet Vega 2300', 'toyota corona', 'plymouth custom', 'audi fox', 'ford torino 500', 'amc matador', 'pontiac catalina brougham', 'dodge d300', 'dodge monaco', 'ford country squire (sw)', 'pontiac safari (sw)', 'amc hornet sportabout (sw)', 'chevrolet Vega (sw)', 'pontiac firebird', 'ford mustang', 'mercury capri 2000', 'opel 1900', 'peugeot 304', 'fiat 124b', 'toyota corolla 1200', 'datsun 1200', 'volkswagen model 111', 'toyota corona mark ii (sw)', 'dodge colt (sw)', 'chevrolet malibu', 'volkswagen 113i deluxe sedan', 'peugeot 504', 'audi 100 ls', 'saab 990', 'bmw 280i', 'amc gremlin', 'ford f250', 'chevy c20', 'dodge d200', 'hi 1200d', 'chevrolet Vega 2300', 'toyota corona', 'plymouth custom', 'audi fox', 'ford torino 500', 'amc matador', 'pontiac catalina brougham', 'dodge d300', 'dodge monaco', 'ford country squire (sw)', 'pontiac safari (sw)', 'amc hornet sportabout (sw)', 'chevrolet Vega (sw)', 'pontiac firebird', 'ford mustang', 'mercury capri 2000', 'opel 1900', 'peugeot 304', 'fiat 124b', 'toyota corolla 1200', 'datsun 1200', 'volkswagen model 111', 'toyota corona mark ii (sw)', 'dodge colt (sw)', 'chevrolet malibu', 'volkswagen 113i deluxe sedan', 'peugeot 504', 'audi 100 ls', 'saab 990', 'bmw 280i', 'amc gremlin', 'ford f250', 'chevy c20', 'dodge d200', 'hi 1200d', 'chevrolet Vega 2300', 'toyota corona', 'plymouth custom', 'audi fox', 'ford torino 500', 'amc matador', 'pontiac catalina brougham', 'dodge d300', 'dodge monaco', 'ford country squire (sw)', 'pontiac safari (sw)', 'amc hornet sportabout (sw)', 'chevrolet Vega (sw)', 'pontiac firebird', 'ford mustang', 'mercury capri 2000', 'opel 1900', 'peugeot 304', 'fiat 124b', 'toyota corolla 1200', 'datsun 1200', 'volkswagen model 111', 'toyota corona mark ii (sw)', 'dodge colt (sw)', 'chevrolet malibu', 'volkswagen 113i deluxe sedan', 'peugeot 504', 'audi 100 ls', 'saab 990', 'bmw 280i', 'amc gremlin', 'ford f250', 'chevy c20', 'dodge d200', 'hi 1200d', 'chevrolet Vega 2300', 'toyota corona', 'plymouth custom', 'audi fox', 'ford torino 500', 'amc matador', 'pontiac catalina brougham', 'dodge d300', 'dodge monaco', 'ford country squire (sw)', 'pontiac safari (sw)', 'amc hornet sportabout (sw)', 'chevrolet Vega (sw)', 'pontiac firebird', 'ford mustang', 'mercury capri 2000', 'opel 1900', 'peugeot 304', 'fiat 124b', 'toyota corolla 1200', 'datsun 1200', 'volkswagen model 111', 'toyota corona mark ii (sw)', 'dodge colt (sw)', 'chevrolet malibu', 'volkswagen 113i deluxe sedan', 'peugeot 504', 'audi 100 ls', 'saab 990', 'bmw 280i', 'amc gremlin', 'ford f250', 'chevy c20', 'dodge d200', 'hi 1200d', 'chevrolet Vega 2300', 'toyota corona', 'plymouth custom', 'audi fox', 'ford torino 500', 'amc matador', 'pontiac catalina brougham', 'dodge d300', 'dodge monaco', 'ford country squire (sw)', 'pontiac safari (sw)', 'amc hornet sportabout (sw)', 'chevrolet Vega (sw)', 'pontiac firebird', 'ford mustang', 'mercury capri 2000', 'opel 1900', 'peugeot 304', 'fiat 124b', 'toyota corolla 1200', 'datsun 1200', 'volkswagen model 111', 'toyota corona mark

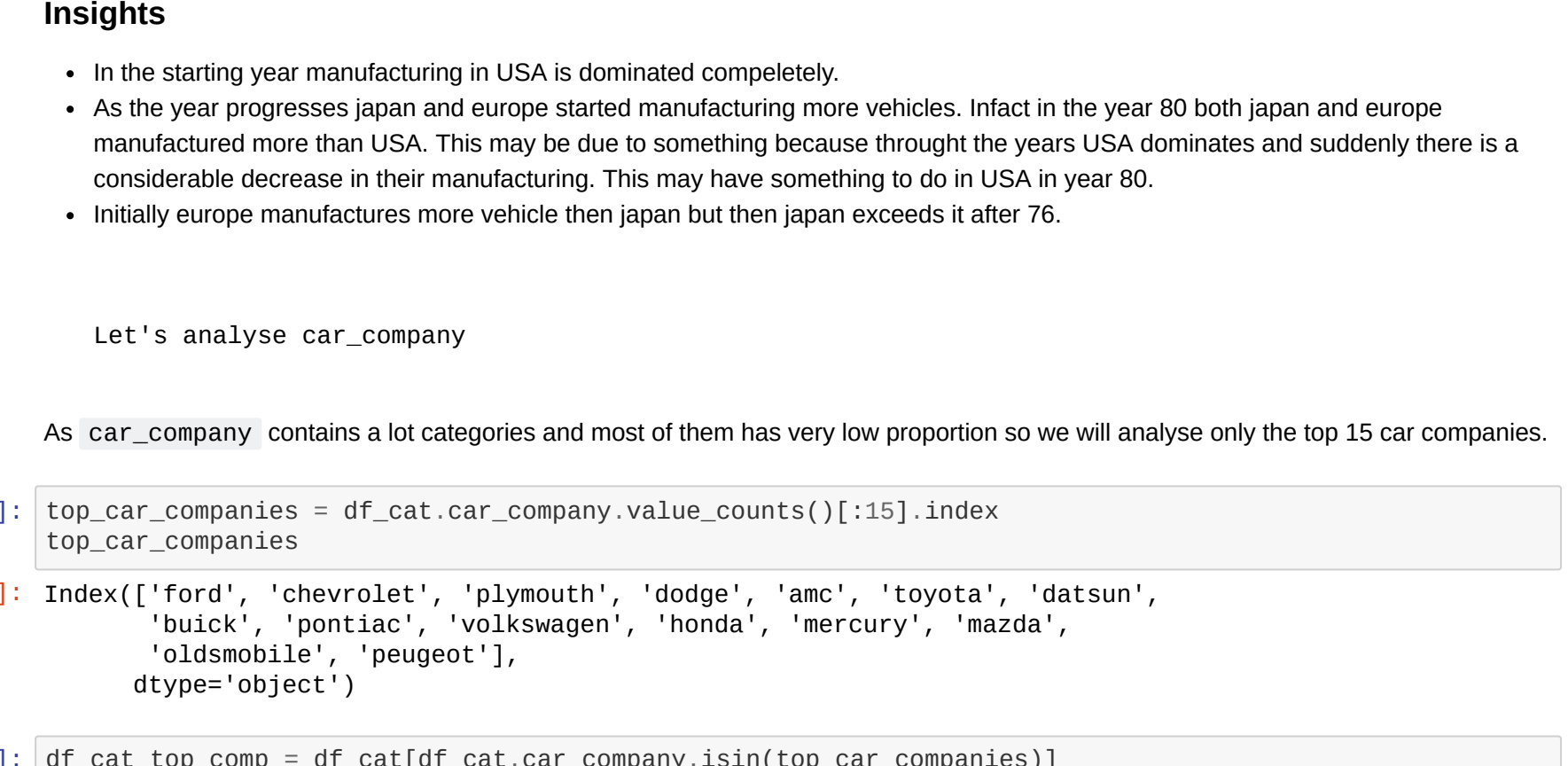

```
[36]: fig = plt.figure(1, (18,4))
sns.countplot(x='model_year', hue='mpg_level', data=df_cat)
plt.show()
```



Insights

- As year progresses manufacturing of low mpg car decreases and eventually stopped after 79.
- As year progresses manufacturing of high mpg car increases and infact after 79 their manufacturing was 1.5-2 times as compared to vehicles with medium mpg.
- Throughout the years there is no significant change in manufacturing of vehicles with medium mpg and infact from 74-79 their manufacturing was more than the sum of remaining two.
- From the scatter plot we can clearly see an upward linear trend i.e., as year progress mpg increases.

```
In [37]: fig = plt.figure(1, (18,4))
sns.countplot(x='model_year', hue='cylinders', data=df_cat)
plt.show()
```

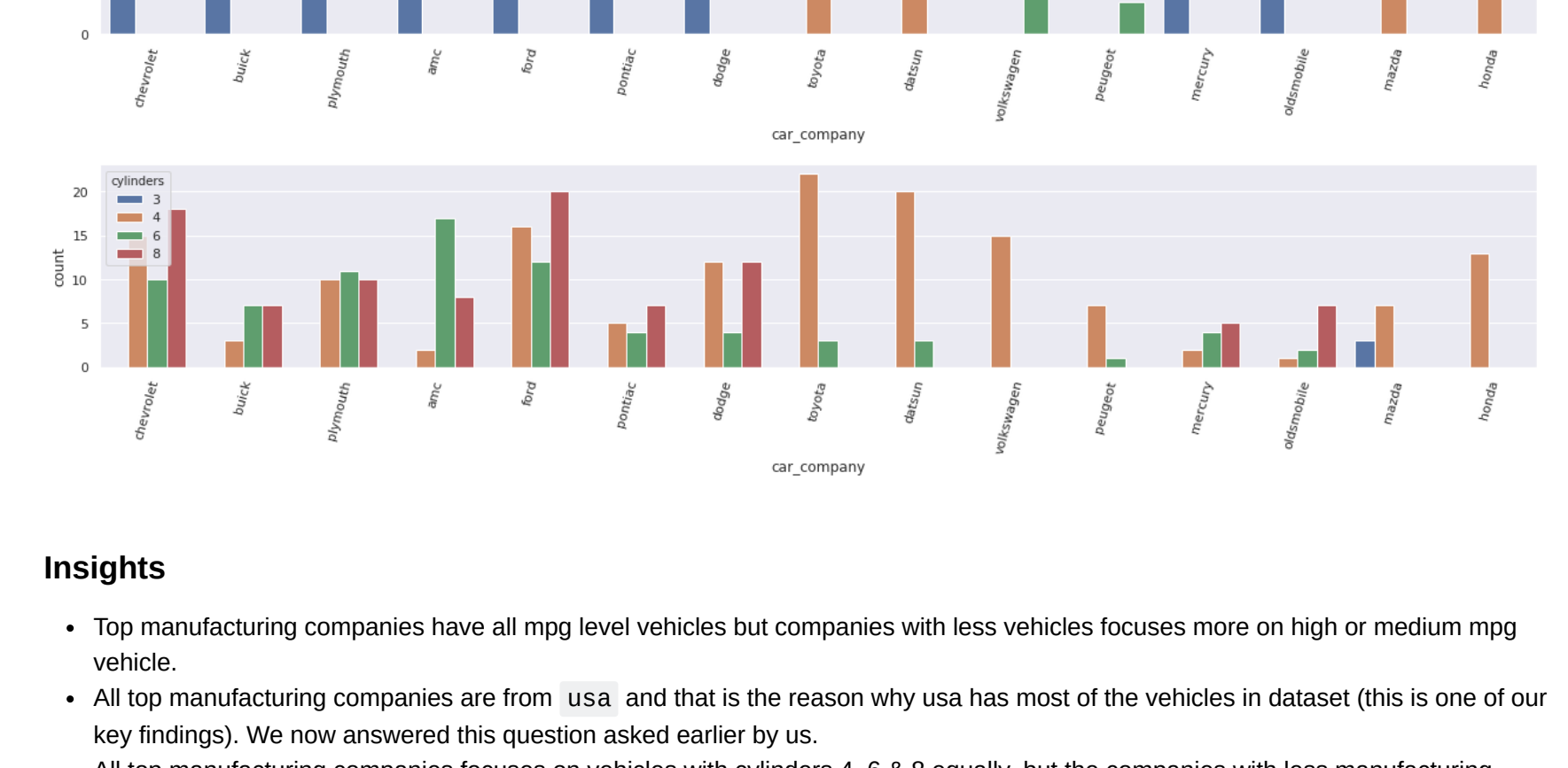


Insights

- As the year progresses vehicles with more cylinders (8 & 6) decreases significantly.
- As the year progresses vehicles with less cylinders increases.
- One important thing to be noticed that throughout the years vehicles with 4 cylinders have significant proportion and infact in the 80's most of the vehicles has 4 cylinders.

These results make sense because as year progresses technology advances vehicles with low mpg and more cylinders loses focus and vehicles with high mpg and less cylinders are the new stars.

```
In [38]: fig = plt.figure(1, (18,4))
sns.countplot(x='model_year', hue='origin', data=df_cat)
plt.show()
```



Insights

- In the starting year manufacturing in USA is dominated completely.
- As the year progresses japan and europe started manufacturing more vehicles. Infact in the year 80 both japan and europe manufactured more than USA. This may be due to something because through the years USA dominates and suddenly there is a considerable decrease in their manufacturing. This may have something to do in USA in year 80.
- Initially europe manufactures more vehicle then japan but then japan exceeds it after 76.

Let's analyse car_company

As car_company contains a lot of categories and most of them have very low proportion so we will analyse only the top 15 car companies.

```
In [39]: top_car_companies = df_cat.car_company.value_counts()[1:15].index
top_car_companies
```

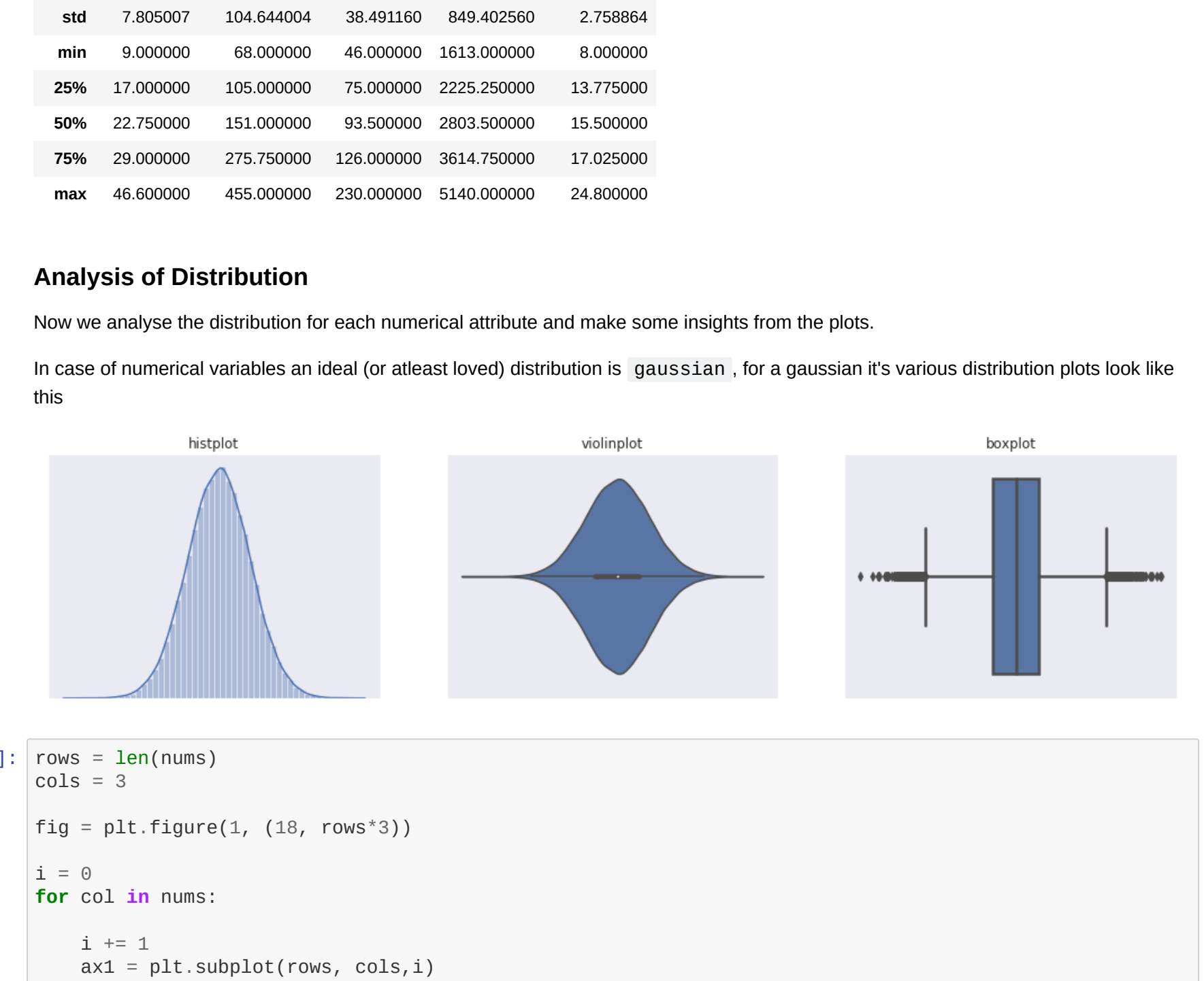
```
Out[39]: Index(['ford', 'chevrolet', 'plymouth', 'dodge', 'amc', 'toyota', 'datsun',
             'buick', 'pontiac', 'volkswagen', 'honda', 'mercury', 'mazda',
             'oldsmobile', 'peugeot'],
            dtype='object')
```

```
In [40]: df_cat_top_comp = df_cat[df_cat.car_company.isin(top_car_companies)]
df_cat_top_comp.shape
```

```
Out[40]: (325, 6)
```

We can see that top 15 car companies alone manufactures 83% of vehicles.

```
In [41]: fig = plt.figure(1, (18,12))
for i,cat in enumerate(['mpg_level', 'origin', 'cylinders']):
    ax = plt.subplot(3,1,i+1)
    sns.countplot(x=cat, hue=cat, data=df_cat_top_comp)
    ax.set_xticklabels(ax.get_xticklabels(), rotation=75)
    plt.tight_layout()
```



Insights

- Top manufacturing companies have all mpg level vehicles but companies with less vehicles focus more on high or medium mpg vehicle.
- All top manufacturing companies are from USA and that is the reason why USA has most of the vehicles in the dataset (this is one of our key findings). We now answered this question asked earlier by us.
- All top manufacturing companies focus on vehicles with cylinders 4, 6 & 8 equally, but the companies with less manufacturing generally use less cylinders in their vehicles.

We are done with the analysis of categorical attributes and found lots of interesting things and answered many unknown questions. Now we will incorporate the required changes of df_cat in to df.

Every attribute except car_name is of interest and participated in our analysis. So we will not add car_name to our dataframe as it is of no interest. This is feature reduction and is an integral part of feature engineering.

```
In [42]: df = pd.concat([df_cat.loc[:, 'origin':'car_company'], df.loc[:, 'mpg':'acceleration']], axis=1)
df.head()
```

```
Out[42]:
```

	origin	cylinders	model_year	mpg_level	car_company	mpg	displacement	horsepower	weight	acceleration
0	usa	8	70	medium	chevrolet	18.0	307.0	130.0	3504	12.0
1	usa	8	70	low	buick	15.0	350.0	165.0	3693	11.5
2	usa	8	70	medium	plymouth	18.0	318.0	150.0	3436	11.0
3	usa	8	70	low	amc	16.0	304.0	150.0	3433	12.0
4	usa	8	70	medium	ford	17.0	302.0	140.0	3449	10.5

```
In [43]: # save these changes to new file
df.to_csv('mpg_cated.csv', index=False)
```

Analysis on Numerical Attributes

Our analysis includes both descriptive stats and EDA.

```
In [44]: df = pd.read_csv('mpg_cated.csv')
df.head()
```

```
Out[44]:
```

	origin	cylinders	model_year	mpg_level	car_company	mpg	displacement	horsepower	weight	acceleration
0	usa	8	70	medium	chevrolet	18.0	307.0	130.0	3504	12.0
1	usa	8	70	low	buick	15.0	350.0	165.0	3693	11.5
2	usa	8	70	medium	plymouth	18.0	318.0	150.0	3436	11.0
3	usa	8	70	low	amc	16.0	304.0	150.0	3433	12.0
4	usa	8	70	medium	ford	17.0	302.0	140.0	3449	10.5

```
In [45]: print('numerical variables: (nums)')
numerical variables: ['mpg', 'displacement', 'horsepower', 'weight', 'acceleration']
```

We will first slice out the numerical columns from original dataframe & then do analysis on it keeping the original data untouched, and at the end incorporate needed changes in our original dataframe.

```
In [46]: df_num = df.loc[:, 'mpg:']
```

Let's first look at some descriptive stats.

```
In [47]: df_num.describe()
```

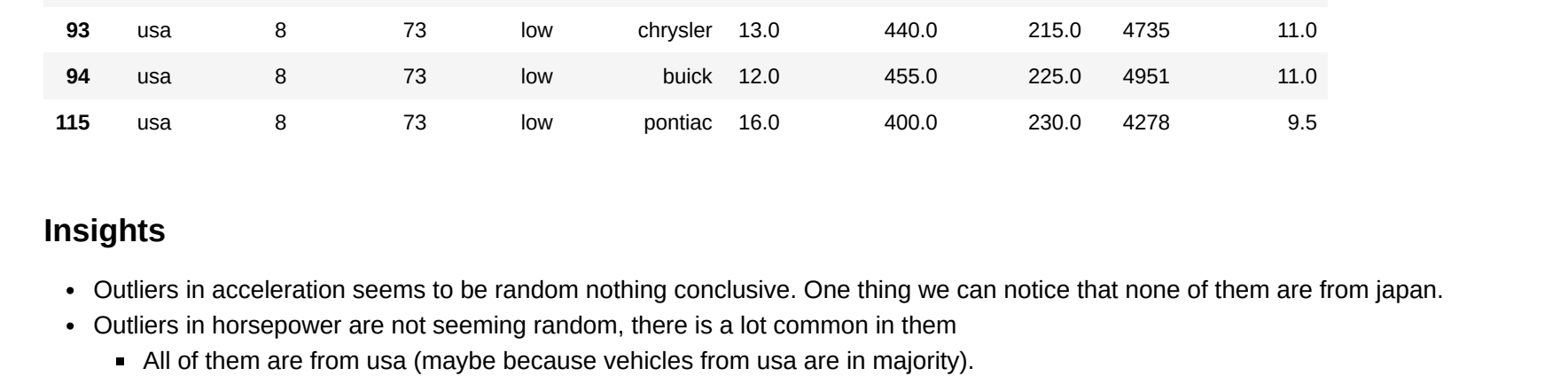
```
Out[47]:
```

	mpg	displacement	horsepower	weight	acceleration
count	392.000000	392.000000	392.000000	392.000000	392.000000
mean	23.449918	194.411990	104.469388	2977.584184	15.541327
std	7.805007	104.644004	38.491160	849.402560	2.758864
min	9.000000	68.000000	46.000000	1613.000000	8.000000
25%	17.000000	105.000000	75.000000	2225.500000	13.775000
50%	22.750000	151.000000	93.500000	2803.500000	15.500000
75%	29.000000	275.750000	126.000000	3614.750000	17.025000
max	46.000000	455.000000	230.000000	5140.000000	24.800000

Analysis of Distribution

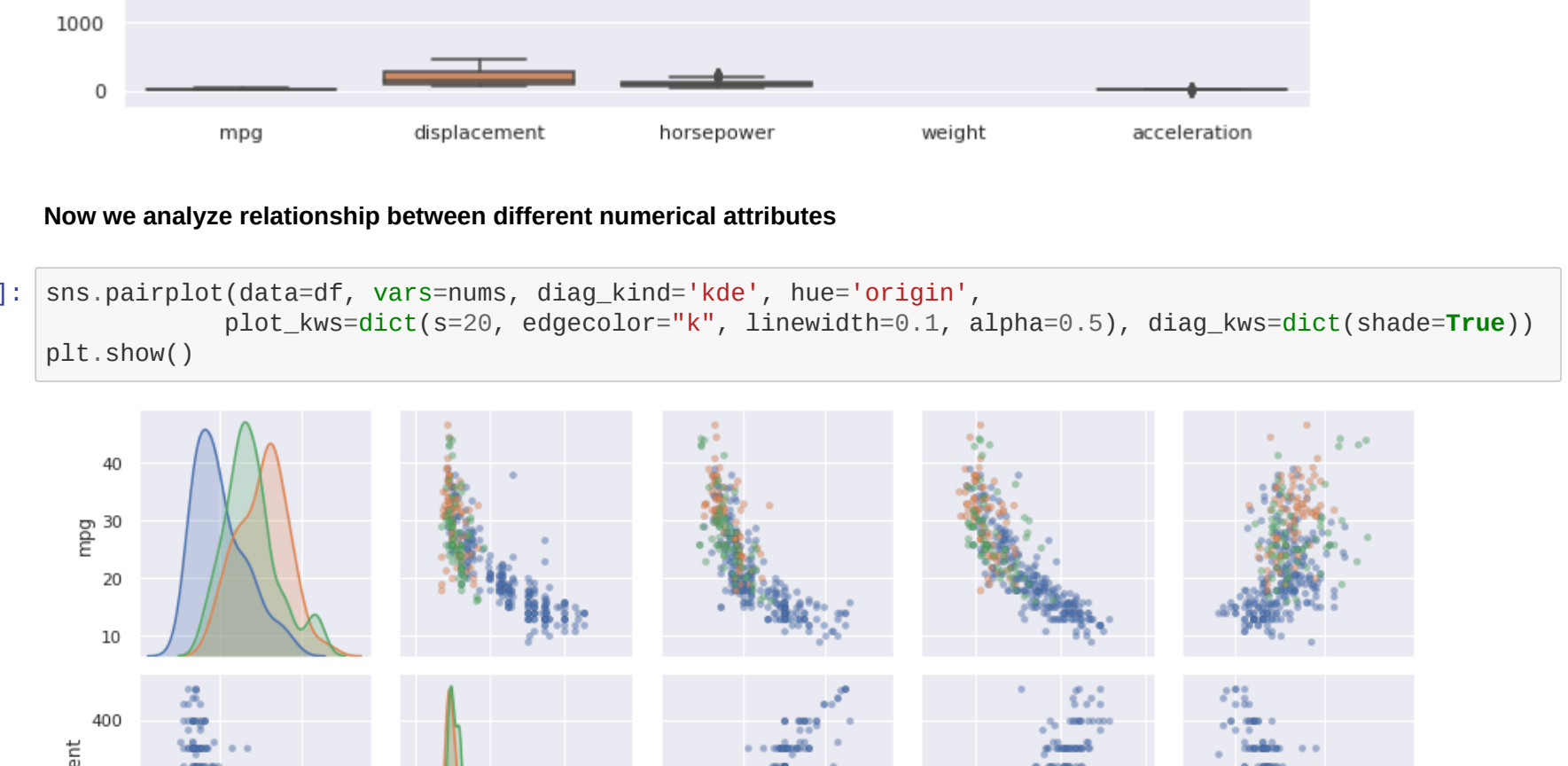
Now we analyse the distribution for each numerical attribute and make some insights from the plots.

In case of numerical variables an ideal (or atleast loved) distribution is gaussian, for a gaussian it's various distribution plots look like this



```
In [48]: rows = len(nums)
cols = 3
```

```
fig = plt.figure(1, (18, rows*3))
i = 0
for col in nums:
    i += 1
    ax1 = plt.subplot(rows, cols, i)
    ax1.hist(df_num[col], alpha=0.6)
    sns.distplot(df_num[col])
    ax1.set_xlabel(None)
    ax1.set_title(f'{col} distribution')
    plt.tight_layout()
    i += 1
    ax2 = plt.subplot(rows, cols, i)
    sns.violinplot(df_num[col])
    sns.swarmplot(df_num[col], alpha=0.6, color='k')
    ax2.set_xlabel(None)
    ax2.set_title(f'{col} swarm-violin plot')
    plt.tight_layout()
    i += 1
    ax3 = plt.subplot(rows, cols, i)
    sns.boxplot(df_num[col], orient='h', linewidth=2.5)
    ax3.set_xlabel(None)
    ax3.set_title(f'{col} box plot')
    plt.tight_layout()
```



Insights

- acceleration is the only distribution which is gaussian. There are few values in acceleration which lie outside the whiskers (the bars extending outwards from the box), these are outliers.
- distributions of mpg & weight seems to be right-skewed gaussian.
- distributions of displacement & horsepower seems to be far from gaussian.

Currently we are analysing the distributions just from plots, in the next phase (statistical analysis) we will do hypothesis testing for the normality of these distributions.

Let's analyse the outliers using tukey formula.

```
In [49]: def tukey_outliers(x):
q1 = np.percentile(x, 25)
q3 = np.percentile(x, 75)
iqr = q3 - q1
min_range = q1 - iqr*1.5
max_range = q3 + iqr*1.5
outliers = x[(x<min_range) | (x>max_range)]
return outliers
```

```
In [50]: for col in nums:
outliers = tukey_outliers(df_num[col])
if len(outliers):
    print(f'{col} has these tukey outliers,\n{outliers}\n')
else:
    print(f'{col} doesn't have any tukey outliers.\n')
```

* displacement doesn't have any tukey outliers.

* horsepower has these tukey outliers,

```
6    220.0
7    215.0
8    225.0
13   225.0
25   215.0
27   218.0
66   288.0
93   215.0
94   225.0
115  230.0
Name: horsepower, dtype: float64
```

* weight doesn't have any tukey outliers.

* acceleration has these tukey outliers,

```
7     8.5
9     8.0
11    8.0
58   23.5
193  22.2
194  22.1
207  21.9
297  24.8
298  22.2
324  23.7
328  24.5
Name: acceleration, dtype: float64
```

acceleration and horsepower are the only attributes with tukey outliers and we can also notice this from the above boxplots.

```
In [51]: df.iloc[list(tukey_outliers(df_num.acceleration).index)]
```

```
Out[51]:
```

	origin	cylinders	model_year	mpg_level	car_company	mpg	displacement	horsepower	weight	acceleration
7	usa	8	70	low	plymouth	14.0	440.0	215.0	4312	8.5
11	usa	8	70	low	amc	15.0	390.0	190.0	3505	8.5
9	usa	8	70	low	plymouth	14.0	340.0	160.0	3609	8.0
193	europe	4	72	medium	volkswagen	28.0	97.0	54.0	2254	22.2
194	usa	4	76	medium	chevrolet	24.5	98.0	60.0	2164	22.1
207	europe	4	76	medium	peugeot	19.0	120.0	88.0	3270	21.9
297	europe	4	79	medium	peugeot	27.2	141.0	71.0	3190	24.8
298	usa	8	79	medium	oldsmobile	23.9	260.0	90.0	3420	22.2
324	europe	4	80	high	vw	43.4	90.0	48.0	2335	23.7
328	europe	4	82	high	vw	44.0	97.0	52.0	2130	24.6

```
In [52]: df.iloc[list(tukey_outliers(df_num.horsepower).index)]
```

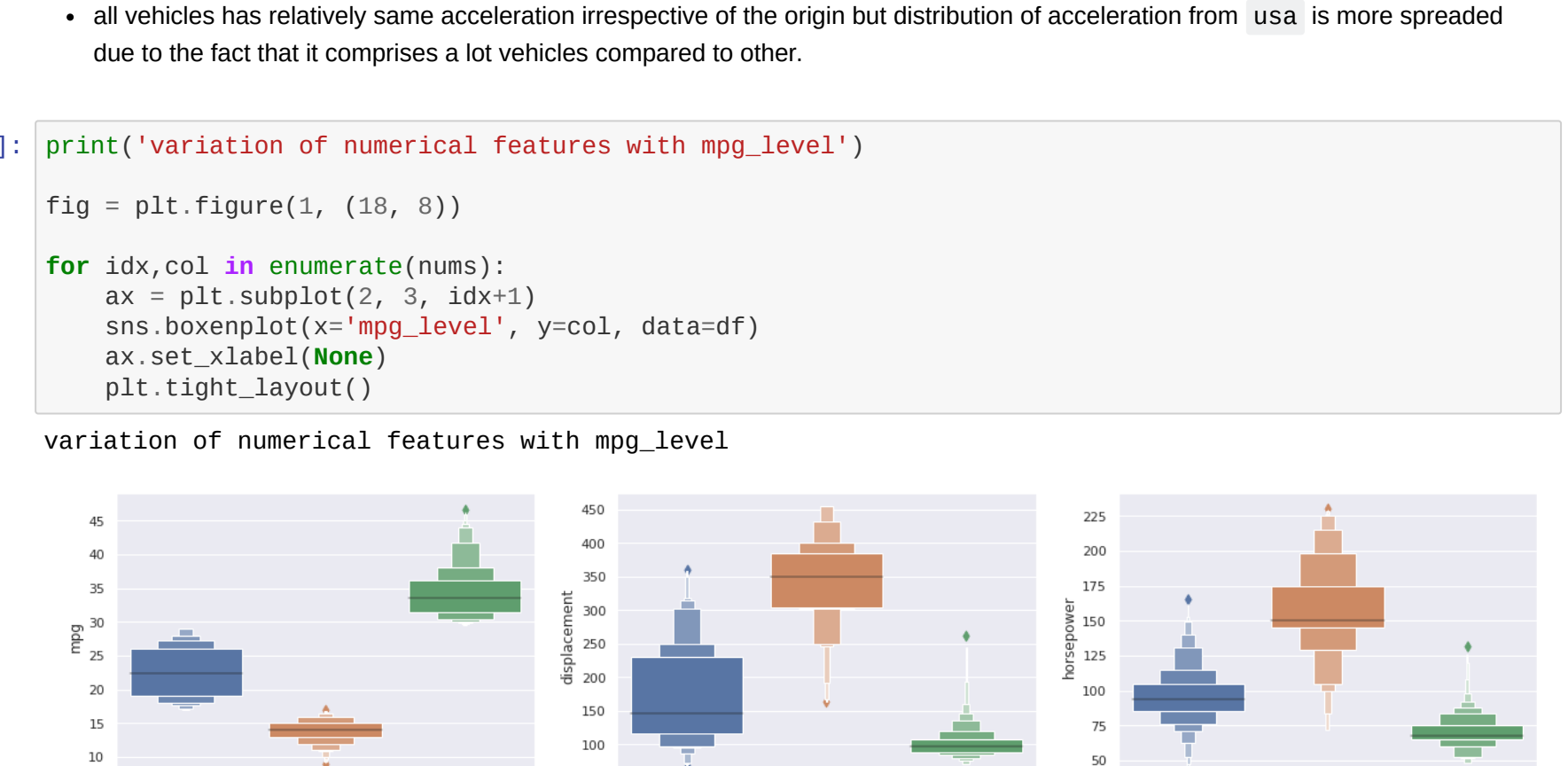
```
Out[52]:
```

	origin	cylinders	model_year	mpg_level	car_company	mpg	displacement	horsepower	weight	acceleration
6	usa	8	70	low	chevrolet	14.0	454.0	220.0	4354	9.0
7	usa	8	70	low	plymouth	14.0	440.0	215.0	4312	8.5
8	usa	8	70	low	portiac	14.0	455.0	225.0	4425	10.0
13	usa	8	70	low	buick	14.0	455.0	230.0	3096	10.0
25	usa	8	70	low	ford	14.0	360.0	215.0	4615	14.0
27	usa	8	70	low	oldodge	11.0	318.0	210.0	4382	13.5
66	usa	8	72	low	mercury	11.0	429.0	208.0	4633	11.0
93	usa	8	73	low	chrysler	13.0	440.0	215.0	4735	11.0
94	usa	8	73	low	buick	12.0	455.0	225.0	4951	11.0
115	usa	8	73	low	portiac	16.0	400.0	230.0	4278	9.5

Insights

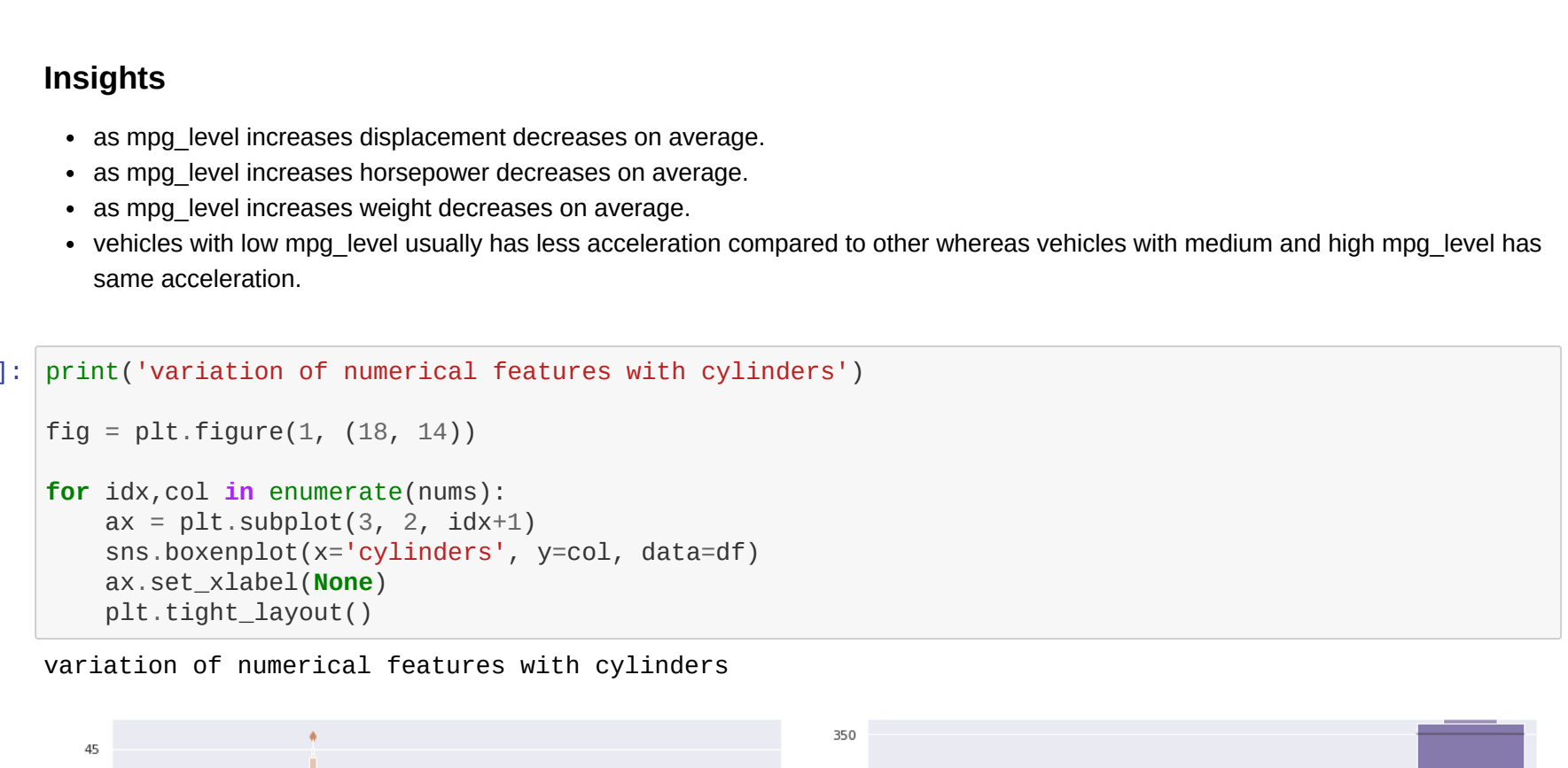
- Outliers in acceleration seems to be random nothing conclusive. One thing we can notice that none of them are from japan.
- Outliers in horsepower are not seeming random, there is a lot common in them
 - All of them are from usa (maybe because vehicles from usa are in majority).
 - All of them has 8 cylinders.
 - All of them has low mpg level.
 - All of them has weight in range 4000.
 - Most of them has displacement in range 400.
 - Most of them were manufactured in early years (before 74).

```
In [53]: # see data is not scaled properly, we need to scale it for modelling but works fine for analysis.
fig = plt.figure(1, (12, 4))
ax = plt.subplot(1, 1, 1)
sns.boxplot(x='variable', y='value', data=pd.melt(df_num))
ax.set_xlabel(None)
ax.set_ylabel(None)
plt.show()
```



Now we analyze relationship between different numerical attributes

```
In [54]: sns.pairplot(data=df, vars=nums, diag_kind='kde', hue='origin',
                 plot_kws=dict(s=20, edgecolor='k', linewidth=0.1, alpha=0.5), diag_kws=dict(shade=True))
plt.show()
```



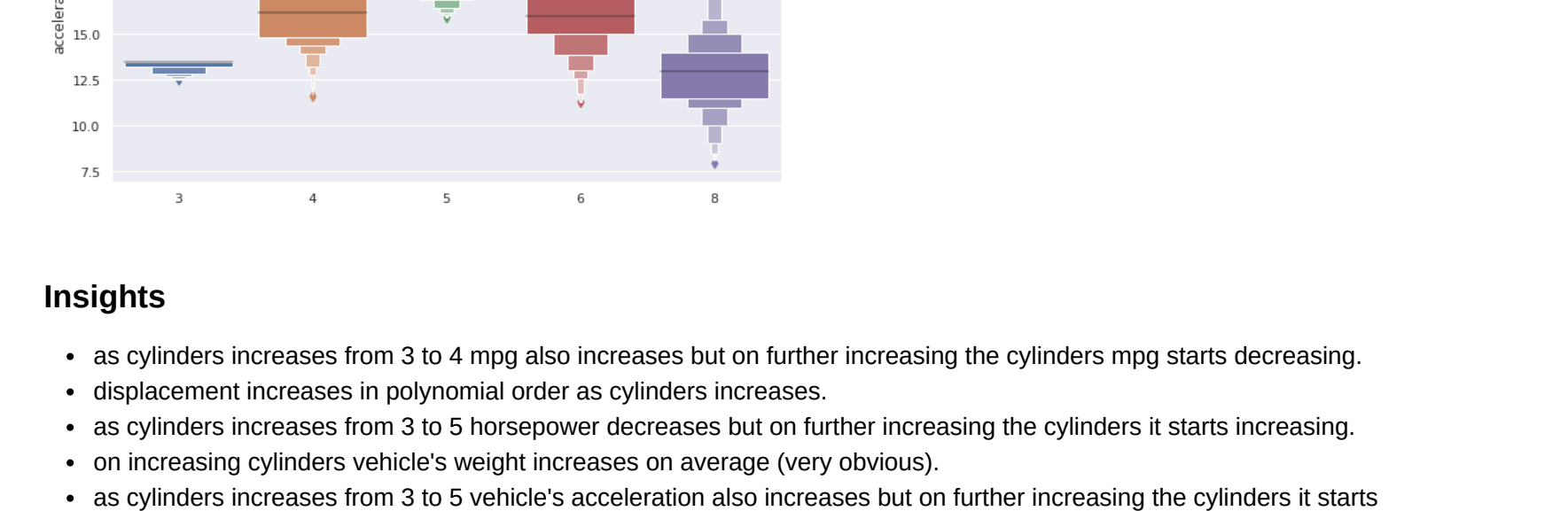
Insights

- As mpg increases displacement, horsepower & weight decreases but acceleration increases.
- As horsepower increases displacement & weight increases but acceleration decreases.
- As weight increases displacement increases but acceleration decreases.
- As acceleration increases displacement decreases.

So all numerical attributes are related with each other.

Now we bifurcate these relationships with different categories.

```
In [56]: '''In this plot we analyze the relationship of horsepower & acceleration
bifurcated by origin, mpg_level & cylinders in a single plot.'''
sns.relplot(x='horsepower', y='acceleration', hue='mpg_level', #style='mpg_level',
            size='cylinders', col='origin', data=df, kind='scatter', sizes=(5, 100), alpha=0.6)
plt.show()
```



Insights

- In every region there is a negative relation b/w horsepower & acceleration.
- vehicles with low mpg has low acceleration and high horsepower.
- vehicles with more cylinders has low acceleration and high horsepower.

```
In [57]: '''In this plot we analyze the relationship of weight & horsepower
bifurcated by origin, mpg_level & cylinders in a single plot.'''
sns.relplot(x='weight', y='horsepower', hue='mpg_level', #style='mpg_level',
            size='cylinders', col='origin', data=df, kind='scatter', sizes=(5, 100), alpha=0.6)
plt.show()
```


Insights

- In every region there is a positive relation b/w weight & horsepower.
- vehicles with low mpg has high weight & horsepower.
- vehicles with more cylinders has high weight & horsepower.

on bifurcating we didn't find anything new.

Now we analyze relationship between numerical and categorical attributes

Note: We are using boxen and violin plots for this but we can also use stripplot.

```
In [58]: print('variation of numerical features with origin')
fig = plt.figure(1, (18, 8))
for idx, col in enumerate(nums):
    ax = plt.subplot(2, 3, idx+1)
    sns.boxenplot(x='origin', y=col, data=df)
    ax.set_xlabel(None)
    plt.tight_layout()
```

variation of numerical features with origin

Insights

- vehicles of usa has less mpg on an average as compared to japan & europe.
- vehicles of usa has more displacement, horsepower and weight as compared to japan & europe.
- All vehicles has relatively same acceleration irrespective of the origin but distribution of acceleration from usa is more spreaded due to the fact that it comprises a lot of vehicles compared to other.

```
In [59]: print('variation of numerical features with mpg_level')
fig = plt.figure(1, (18, 8))
for idx, col in enumerate(nums):
    ax = plt.subplot(2, 3, idx+1)
    sns.boxenplot(x='mpg_level', y=col, data=df)
    ax.set_xlabel(None)
    plt.tight_layout()
```

variation of numerical features with mpg_level

Insights

- As mpg_level increases displacement decreases on average.
- As mpg_level increases horsepower decreases on average.
- As the year progresses mpg_level slowly increases acceleration compared to other whereas vehicles with medium and high mpg_level has same acceleration.

```
In [60]: print('variation of numerical features with cylinders')
fig = plt.figure(1, (18, 14))
for idx, col in enumerate(nums):
    ax = plt.subplot(3, 2, idx+1)
    sns.boxenplot(x='cylinders', y=col, data=df)
    ax.set_xlabel(None)
    plt.tight_layout()
```

variation of numerical features with cylinders

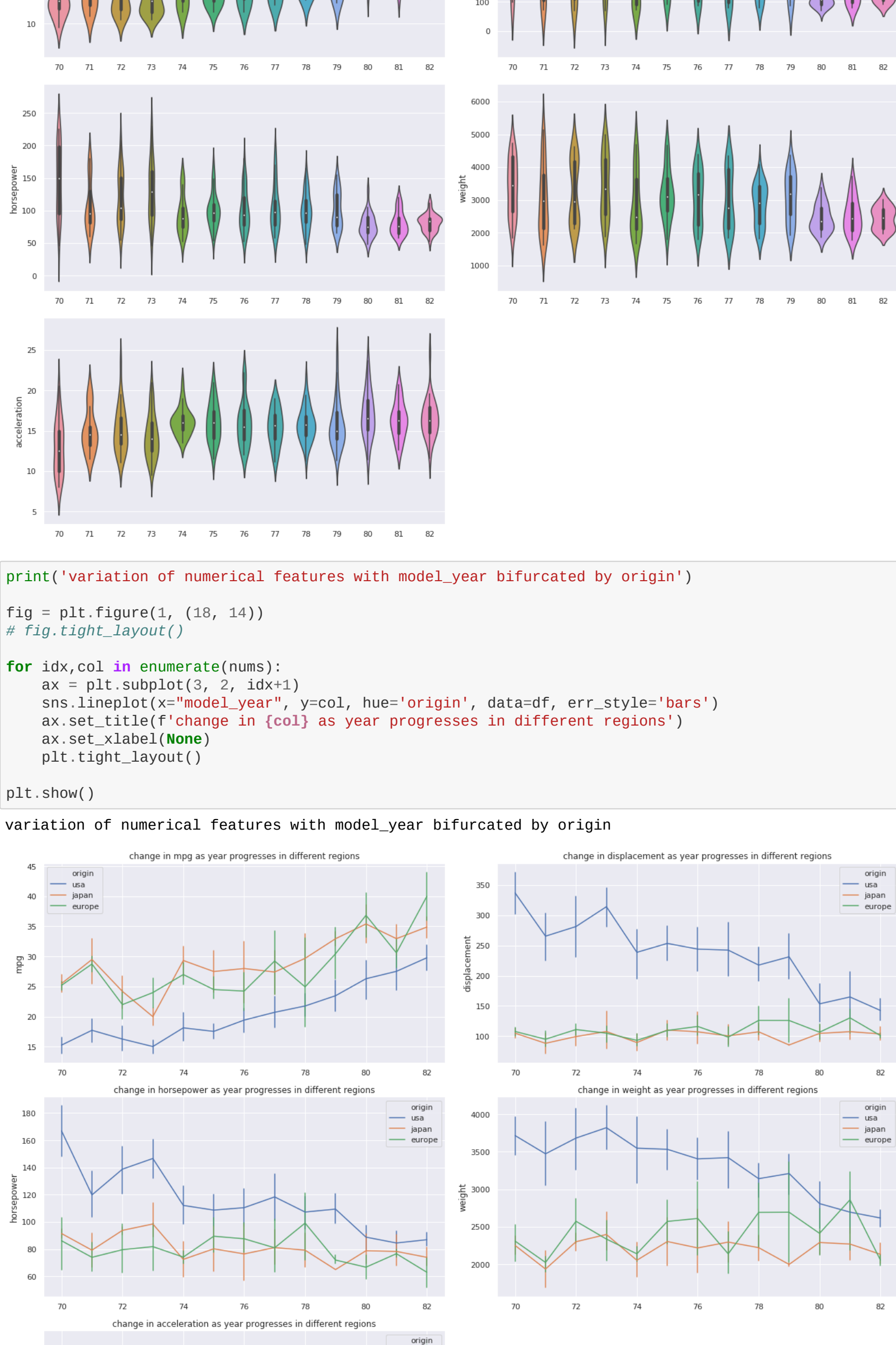
Insights

- As cylinders increases from 3 to 4 mpg also increases but on further increasing the cylinders mpg starts decreasing.
- displacement increases in polynomial order as cylinders increases.
- As cylinders increases from 3 to 5 horsepower decreases but on further increasing the cylinders it starts increasing.
- on increasing cylinders vehicle's weight increases on average (very obvious).
- As cylinders increases from 3 to 5 vehicle's acceleration also increases but on further increasing the cylinders it starts decreasing (maybe due to the fact that vehicles with more cylinders have more weight and hence less acceleration).


```
In [61]: print('variation of numerical features with model_year')
fig = plt.figure(1, (18, 14))
# Fig.tight_layout()
```

```
for idx,col in enumerate(nums):
    ax = plt.subplot(3, 2, idx+1)
    sns.violinplot(x='model_year', y=col, data=df)
    ax.set_xlabel(None)
    plt.tight_layout()
```

variation of numerical features with model_year

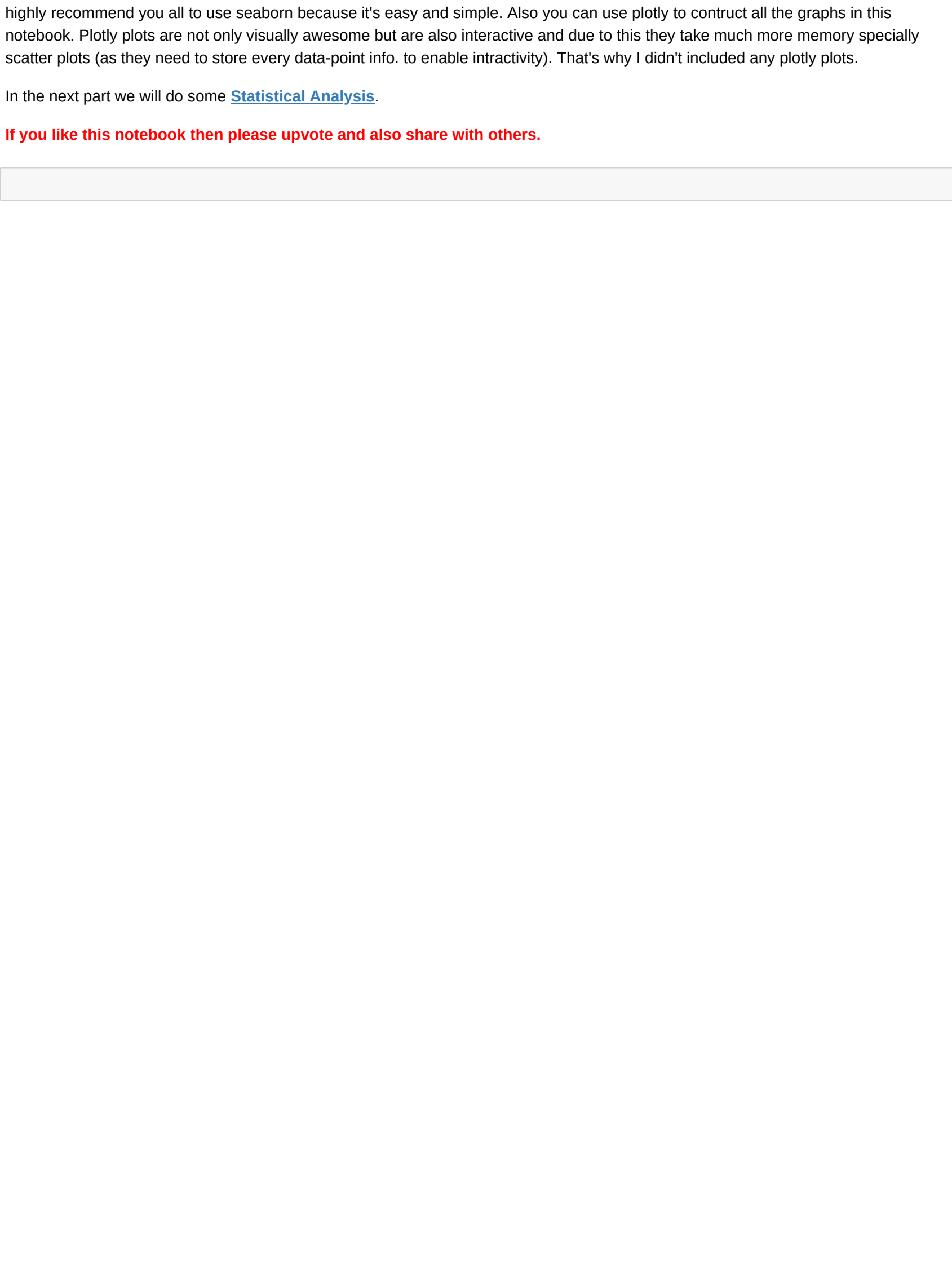


```
In [62]: print('variation of numerical features with model_year bifurcated by origin')
fig = plt.figure(1, (18, 14))
# Fig.tight_layout()

for idx,col in enumerate(nums):
    ax = plt.subplot(3, 2, idx+1)
    sns.lineplot(x='model_year', y=col, hue='origin', data=df, err_style='bars')
    ax.set_title(f'change in {col} as year progresses in different regions')
    ax.set_xlabel(None)
    plt.tight_layout()

plt.show()
```

variation of numerical features with model_year bifurcated by origin



Insights

- as year progresses there is an increase in mpg across all origins (this we already observed in analysis on categorical data).
- as year progresses there is a slight decrease in displacement, horsepower & weight of the vehicles belonging to usa but there is no significant change in japan & europe. One thing we can observe is that in the 80's all vehicles have similar displacement because unlike the 70's the distribution is not spread out (i.e., distribution is **short fatty** instead of **tall skinny**).
- throughout the years acceleration remains relatively the same across all regions.

So we are done for now. We did some good amount of EDA and also explored various plotting features provided to us by seaborn. I highly recommend you all to use seaborn because it's easy and simple. Also you can use `plotly` to construct all the graphs in this notebook. `Plotly` plots are not only visually awesome but are also interactive and due to this they take much more memory specially scatter plots (as they need to store every data-point info. to enable interactivity). That's why I didn't included any `plotly` plots.

In the next part we will do some [Statistical Analysis](#).

If you like this notebook then please upvote and also share with others.

```
In [ ]:
```