



# Matplotlib and Seaborn



# Library for data analysis

---

- **NumPy** (Numerical Python)
  - Large multidimensional array operations
- **SciPy** (Scientific Python)
  - Many efficient numerical routines such as routines for numerical integration and optimization
- **Pandas**
  - Data manipulation and data visualization
- **Matplotlib**
  - Data exploration and data visualization
- **Seaborn**
  - High-level data visualization library based on Matplotlib
- **Scikit-learn**
  - Machine learning and statistical modeling

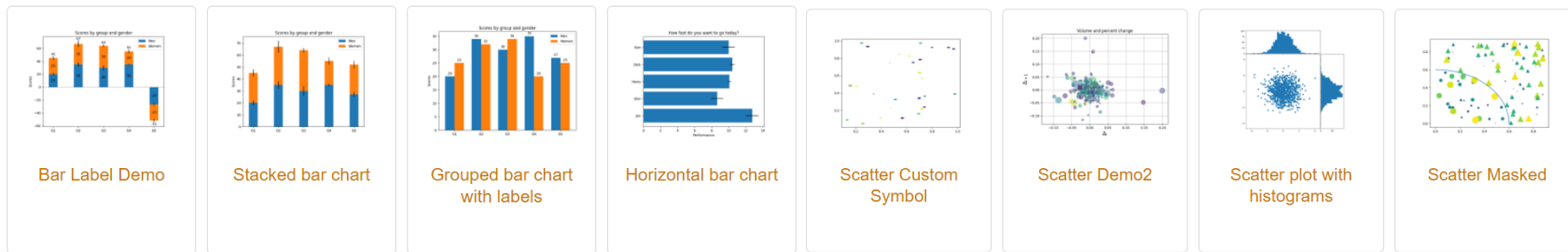
# Outline

---

- Matplotlib
  - Single plot
  - Multiple plot
  - Secondary y-axis
- Seaborn
  - X-axis with categorical data
    - Countplot, barplot, heatmap
  - Numerical data
    - Histogram, scatter plot, joint plot, pair plot
  - Facet Grid

# Matplotlib

- **matplotlib**
  - Matplotlib is a **low-level** data visualization library in python.

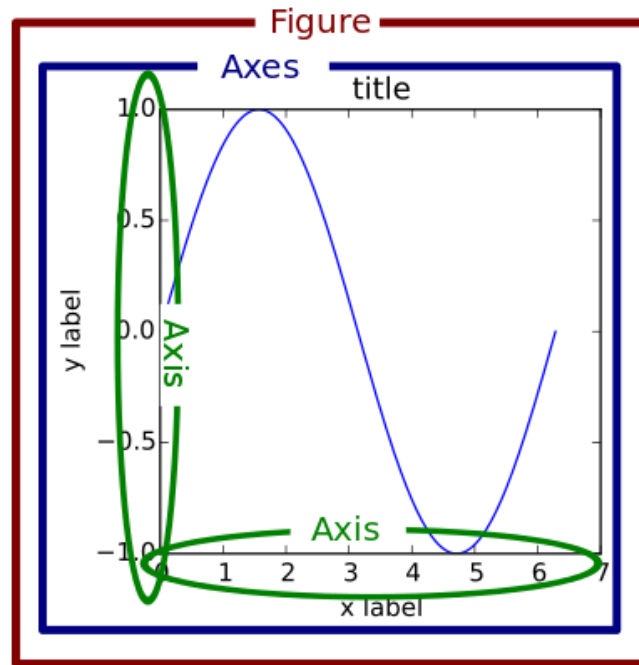


- **matplotlib.pyplot**
  - Pyplot is a subset of matplotlib, which is a collection of the most commonly used plotting functions.

```
import matplotlib.pyplot as plt
```

# Matplotlib - figure and axes

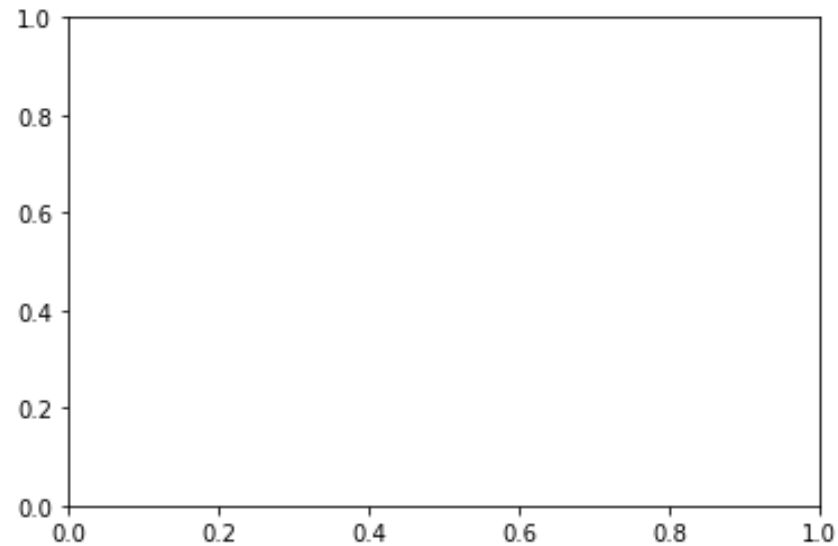
- **Figure** object: The outermost container for matplotlib plots, which can contain multiple Axes objects.
- **Axes** object: The axes is the area your plot appears in. (In matplotlib, axes is not the plural form of axis)



# Matplotlib - single plot (1 / 3)

- Step1: Create figure and axes
- Step2: Plot a chart in axes
- Step3: Format the style

```
#step1: Create a figure and axes  
fig = plt.figure(figsize=(6, 4))  
ax1 = fig.add_subplot()
```



```
type(fig)
```

```
matplotlib.figure.Figure
```

```
type(ax1)
```

```
matplotlib.axes._subplots.AxesSubplot
```

# Matplotlib - single plot (2/3)

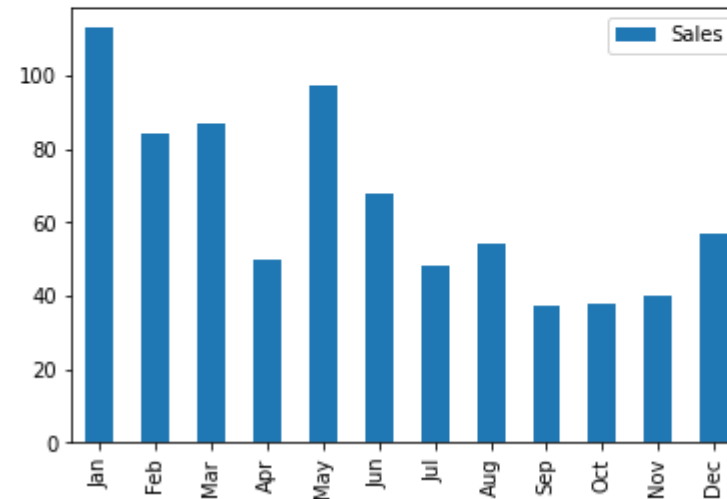
- Step1: Create figure and axes
- Step2: Plot a chart in axes
- Step3: Format the style

	Sales	Cumulative_sales
Jan	113	113
Feb	84	197
Mar	87	284
Apr	50	334
May	97	431
Jun	68	499
Jul	48	547
Aug	54	601
Sep	37	638
Oct	38	676
Nov	40	716
Dec	57	773

```
#step1: Create a figure and axes  
fig = plt.figure(figsize=(6, 4))  
ax1 = fig.add_subplot()
```

```
#step2: Plot a chart in axes  
sales_df.plot(kind = "bar", y = "Sales", ax = ax1)
```

<AxesSubplot:>



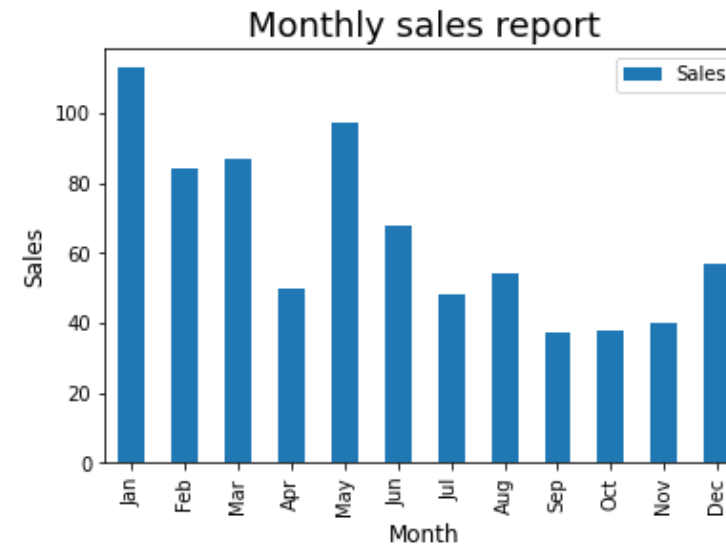
# Matplotlib - single plot (3/3)

- Step1: Create figure and axes
- Step2: Plot a chart in axes
- Step3: Format the style

```
#step1: Create a figure and axes
fig = plt.figure(figsize=(6, 4))
ax1 = fig.add_subplot()

#step2: Plot a chart in axes
sales_df.plot(kind = "bar", y = "Sales", ax = ax1)

#step3: Format the style
ax1.set_title("Monthly sales report", fontsize=18)
ax1.set_xlabel("Month", fontsize=12)
ax1.set_ylabel("Sales", fontsize=12);
```

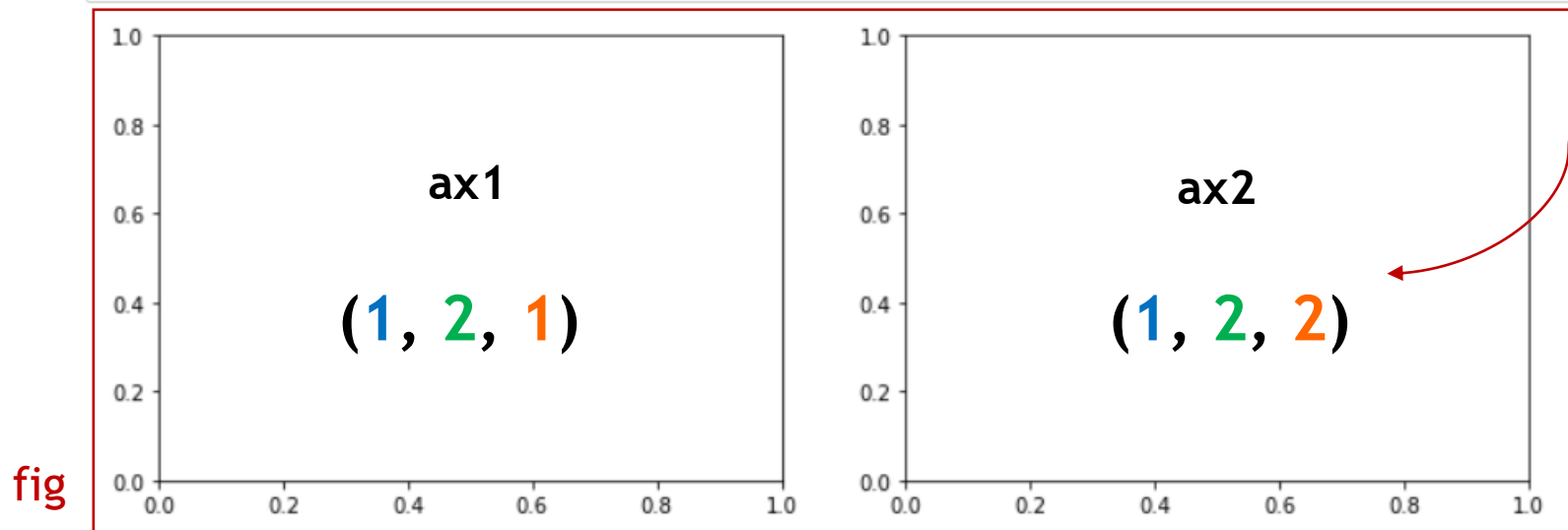




# Matplotlib - multiple plots

- Use three integers (**nrows**, **ncols**, **index**) to describe the positions.
- The subplot will take the **index** position on a grid with **nrows** rows and **ncols** columns. index starts at 1 in the upper left corner and increases to the right.

```
#step1: Create a figure and axes  
fig = plt.figure(figsize=(12,4))  
ax1 = fig.add_subplot(1,2,1)  
ax2 = fig.add_subplot(1,2,2)
```



1 x 2 grid

1 (1,2,1)	2 (1,2,2)
--------------	--------------

2 x 3 grid

1 (2,3,1)	2 (2,3,2)	3 (2,3,3)
4 (2,3,4)	5 (2,3,5)	6 (2,3,6)

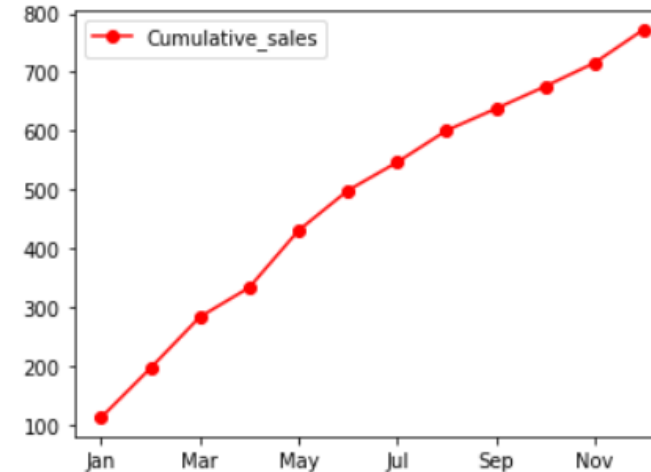
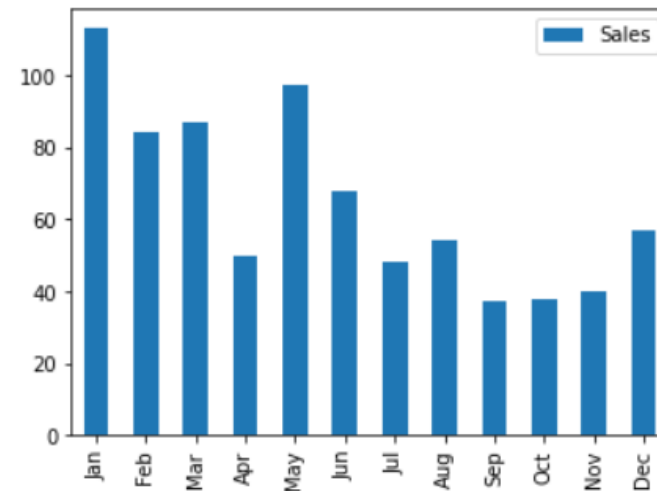
# Matplotlib - multiple plots

- Assign axes by the argument `ax`.

	Sales	Cumulative_sales
Jan	113	113
Feb	84	197
Mar	87	284
Apr	50	334
May	97	431
Jun	68	499
Jul	48	547
Aug	54	601
Sep	37	638
Oct	38	676
Nov	40	716
Dec	57	773

```
#step1: Create a figure and axes
fig = plt.figure(figsize=(12,4))
ax1 = fig.add_subplot(1,2,1)
ax2 = fig.add_subplot(1,2,2)
#step2: Plot a chart in axes
sales_df.plot(kind = "bar", y = "Sales", ax = ax1)
sales_df.plot(kind = "line", y = "Cumulative_sales", ax = ax2, color = "red", marker = "o")
```

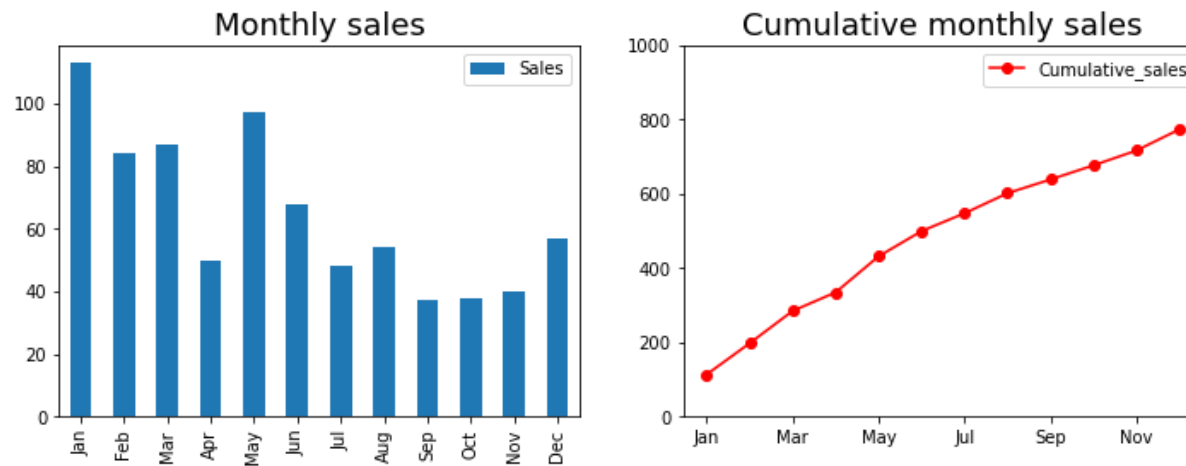
<AxesSubplot:>



# Matplotlib - multiple plots

- Plot charts in axes and format the charts

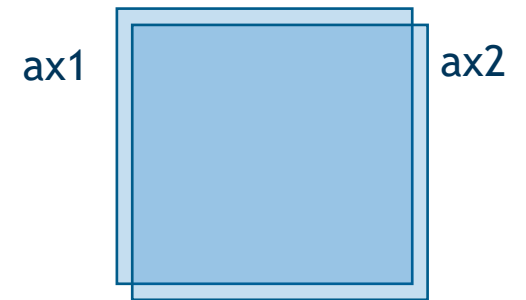
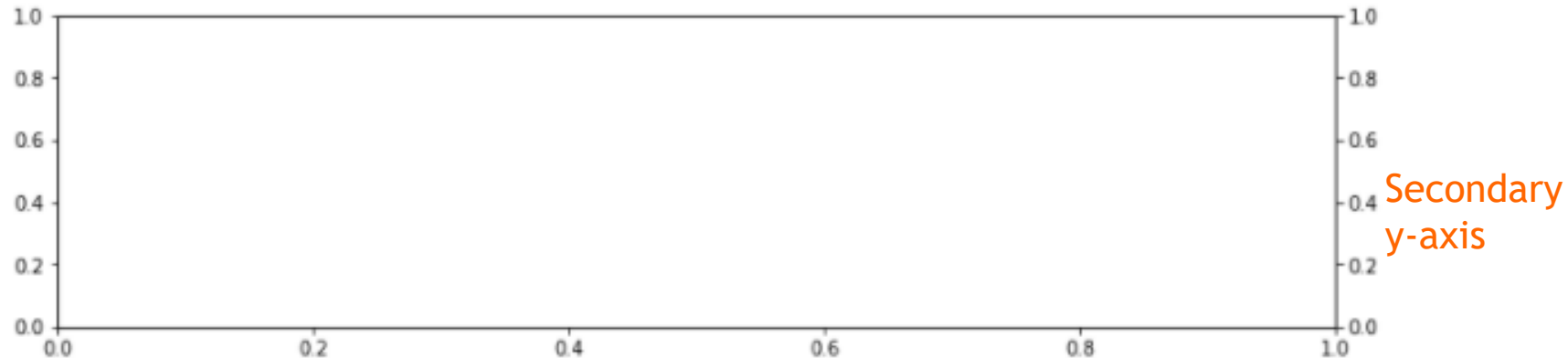
```
#step1: Create a figure and axes
fig = plt.figure(figsize=(12,4))
ax1 = fig.add_subplot(1,2,1)
ax2 = fig.add_subplot(1,2,2)
#step2: Plot a chart in axes
sales_df.plot(kind = "bar", y = "Sales", ax=ax1)
sales_df.plot(kind = "line", y = "Cumulative_sales", ax=ax2, color = "red", marker = "o")
#step3: Format the style
ax1.set_title("Monthly sales", fontsize = 18)
ax2.set_title("Cumulative monthly sales", fontsize = 18)
ax2.set_ylim([0,1000]);
```



# Matplotlib - Single plot with secondary y-axis

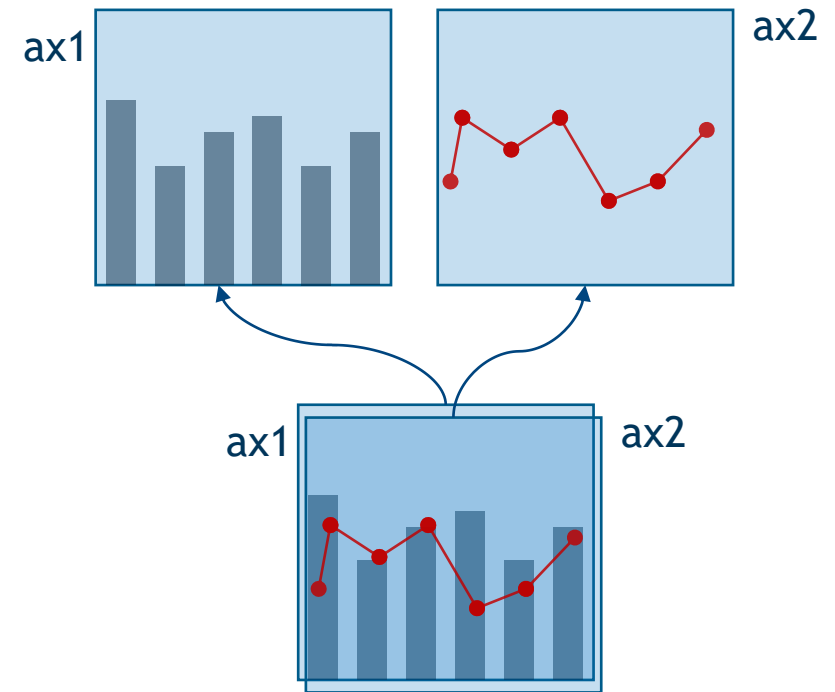
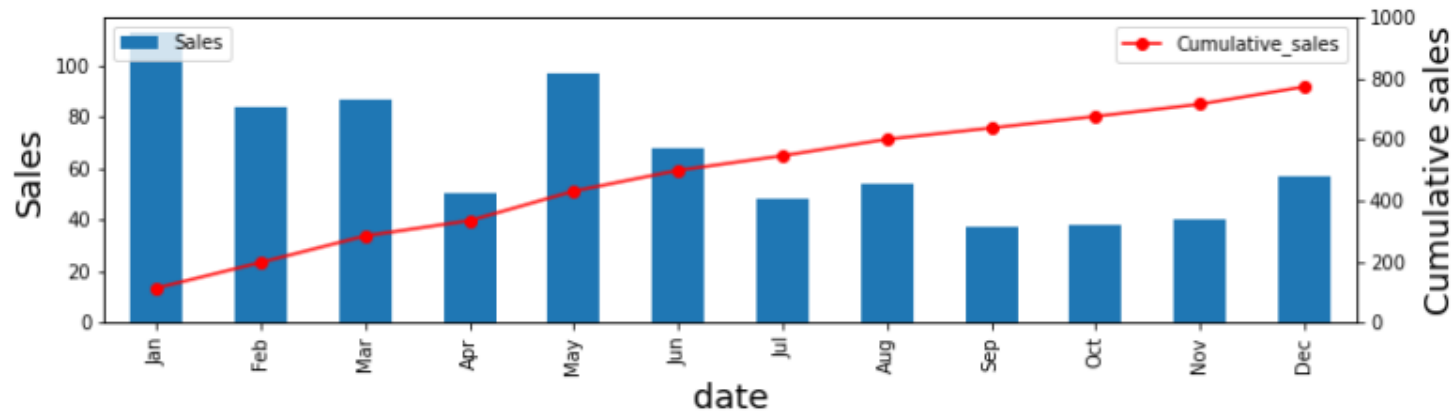
- In some cases, we need the **secondary y-axis**: it allows you to use the same X axis with two different sets of Y-axis data with **two different scales**.

```
#step1: Create a figure and axes  
fig = plt.figure(figsize=(12,3))  
ax1 = fig.add_subplot()  
ax2 = ax1.twinx() Share the same x-axis as ax1
```



# Matplotlib - Single plot with secondary y-axis

```
#step1: Create a figure and axes
fig = plt.figure(figsize=(12,3))
ax1 = fig.add_subplot()
ax2 = ax1.twinx()
#step2: Plot a chart in axes
sales_df.plot(kind = "bar", y = "Sales", ax=ax1)
sales_df.plot(kind = "line", y = "Cumulative_sales", ax=ax2, color = "red", marker = "o")
#step3: Format the style
ax1.set_xlabel("date", fontsize = 18)
ax1.set_ylabel("Sales", fontsize = 18)
ax1.legend(loc="upper left")
ax2.set_ylabel("Cumulative sales", fontsize = 18)
ax2.set_ylim([0,1000])
ax2.legend(loc="upper right");
```



# Exercise

(A.1) Given the synthetic dataset above, each column represents the quarterly sales of products A, B, and C. Create a figure with three bar charts to show the sales data of each product.

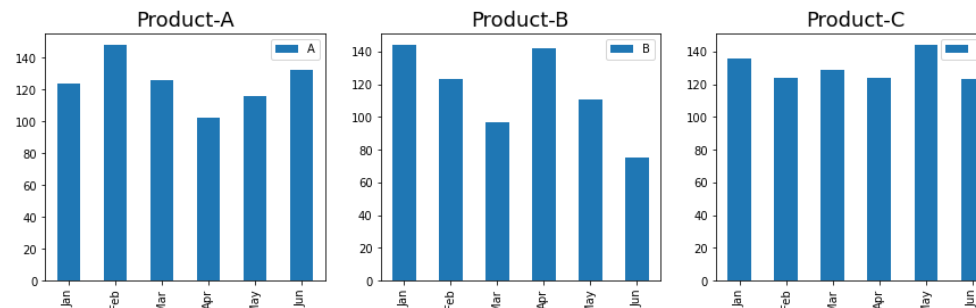
Setting: `figsize = (15,4)`

*#step1: Create a figure and axes*

*#step2: Plot a chart in axes*

*#step3: Format the style*

	A	B	C
Jan	124	144	136
Feb	148	123	124
Mar	126	97	129
Apr	102	142	124
May	116	111	144
Jun	132	75	123



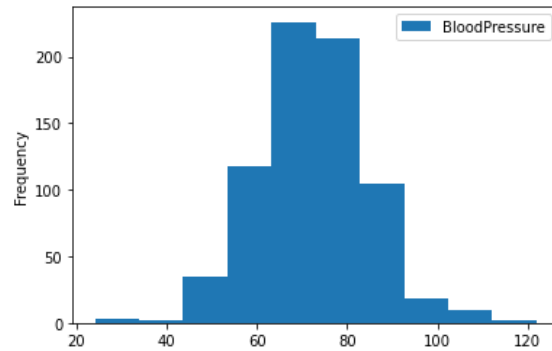
1 x 3 grid

1	2	3
(1,3,1)	(1,3,2)	(1,3,3)

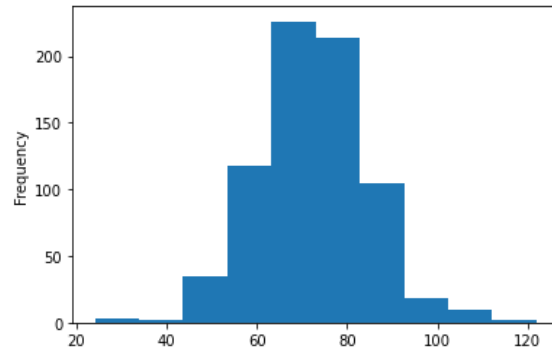
# Different plotting methods

## Pandas

```
diabetes_df.plot(kind = "hist", y = "BloodPressure");
```



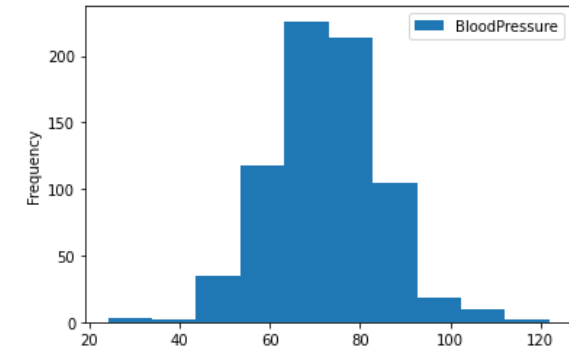
```
plt.hist(diabetes_df.BloodPressure)  
plt.ylabel('Frequency')  
plt.show()
```



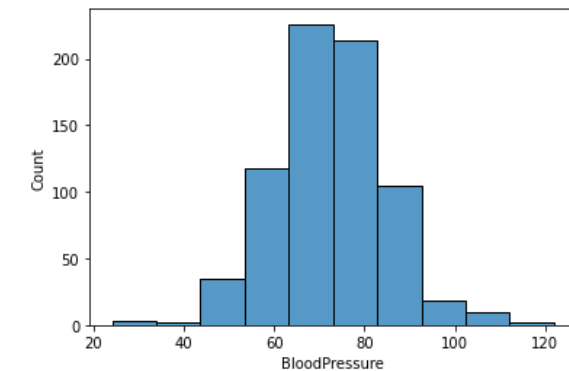
## Matplotlib

## Matplotlib (Create a container) Pandas (plot)

```
fig = plt.figure(figsize=(6, 4))  
ax1 = fig.add_subplot()  
diabetes_df.plot(kind = "hist", y = "BloodPressure", ax=ax1);
```



```
sns.histplot(data = diabetes_df, x = "BloodPressure", bins = 10);
```



# Outline

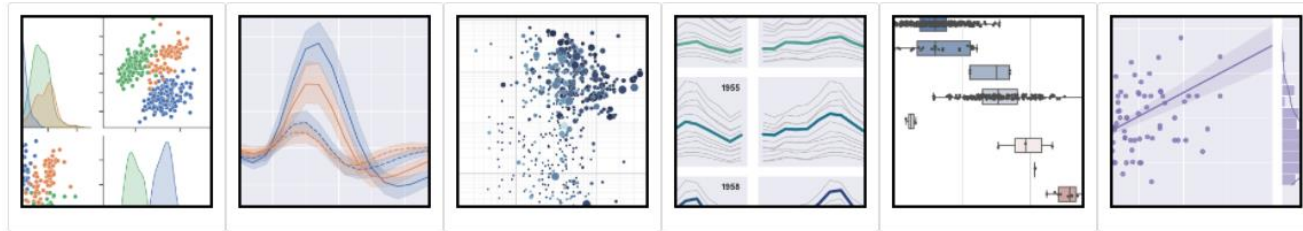
---

- Matplotlib
  - Single plot
  - Multiple plot
  - Secondary y-axis
- Seaborn
  - X-axis with categorical data
    - Countplot, barplot, heatmap
  - Numerical data
    - Histogram, scatter plot, joint plot, pair plot
  - Facet Grid



# Seaborn

- Seaborn is a package that provides many types of insightful plots. You can write simple code to visualize complex graphics.



- Data structures accepted by seaborn
  - Objects from pandas or numpy
  - Long-form and wide-form
- Installation

```
pip install seaborn
```

```
import seaborn as sns
```

# Countplot

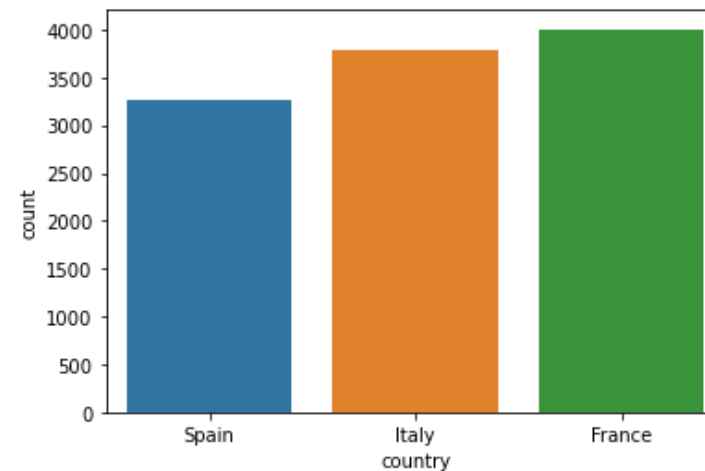
- Use `countplot()` to show the counts of observations in each bar.
- Pass arguments: `data`, `x`

	country	description	designation	points	price	province	region_1	region_2	variety	winery
51	France	This structured, complex Chardonnay is packed...	NaN	90	68.0	Burgundy	Chassagne-Montrachet	NaN	Chardonnay	Chartron et Trébouchet
53	France	With its light color and cool feel, this well...	L'Inédit	90	28.0	Loire Valley	Coteaux du Giennois	NaN	Pinot Noir	Clement et Florian Berthier
63	France	L'Homme Mort is a northern extension of the Fo...	L'Homme Mort Premier Cru	91	45.0	Burgundy	Chablis	NaN	Chardonnay	Domaine Chenevieres
66	France	The steely character of a young Chablis is ver...	Fourchaume Premier Cru	91	38.0	Burgundy	Chablis	NaN	Chardonnay	Louis Max
76	France	This wine is bone-dry, although with some age ...	Le Nombre d'Or Brut Nature	91	85.0	Champagne	Champagne	NaN	Chardonnay	Aubry
...	...	...	...	...	...	...	...	...	...	...
150908	France	Another premier cru from Michel Gros, this one...	Aux Brulees	90	65.0	Burgundy	Vosne-Romanée	NaN	Pinot Noir	Michel Gros
150909	France	This is a lovely, fragrant Burgundy, with a sm...	Clos dea Argillieres	89	52.0	Burgundy	Nuits-St-Georges	NaN	Pinot Noir	Daniel Rion
150910	France	Scents of graham cracker and malted milk choco...	NaN	89	38.0	Burgundy	Chambolle-Musigny	NaN	Pinot Noir	Michel Gros
150911	France	This needs a good bit of breathing time, then ...	Les Chaliots	87	37.0	Burgundy	Nuits-St-Georges	NaN	Pinot Noir	Michel Gros
150912	France	The nose is dominated by the attractive scents...	Les Charmes	87	65.0	Burgundy	Chambolle-Musigny	NaN	Pinot Noir	Daniel Rion

4580 rows x 10 columns

```
sns.countplot(data = wine_df, x = "country")
```

```
<AxesSubplot:xlabel='country', ylabel='count'>
```



The bar height is the number of wines of each country.

No need to group data in advance

```
wine_df.groupby("country").size()
```

```
country
France    4004
Italy     3783
Spain     3258
dtype: int64
```

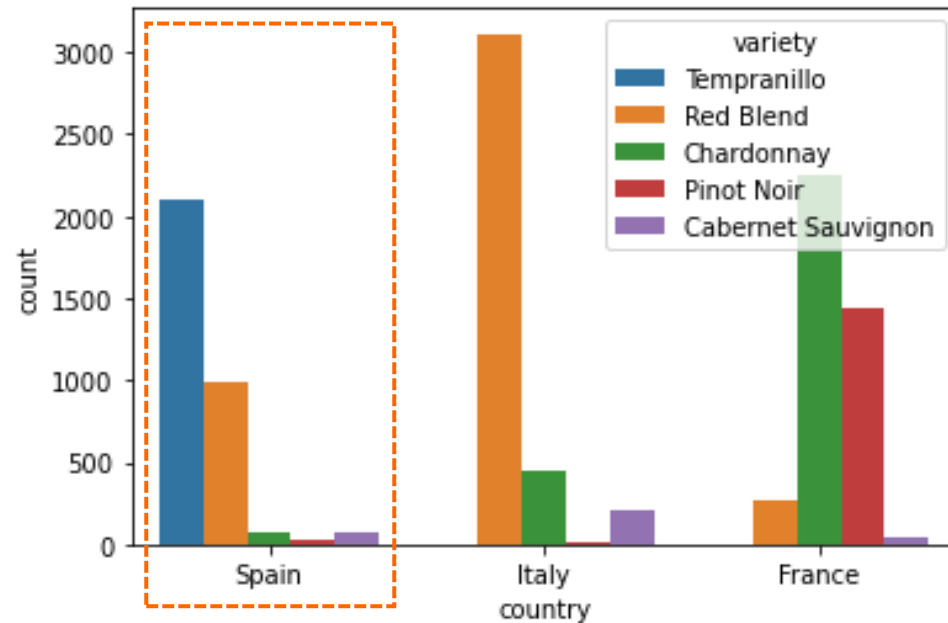
Understand the data behind the chart

# Countplot

- Use the argument `hue` to specify the categorical variable used for color encoding.

```
sns.countplot(data = wine_df, x = "country", hue = "variety")
```

```
<AxesSubplot:xlabel='country', ylabel='count'>
```



```
wine_df.groupby(["country", "variety"]).size()
```

country	variety	count
France	Cabernet Sauvignon	44
	Chardonnay	2247
	Pinot Noir	1446
	Red Blend	267
	Cabernet Sauvignon	214
Italy	Chardonnay	450
	Pinot Noir	7
	Red Blend	3111
	Tempranillo	1
	Spain	Cabernet Sauvignon
Spain	Chardonnay	67
	Pinot Noir	31
	Red Blend	987
	Tempranillo	2899

dtype: int64

Understand the data  
behind the chart

# Barplot

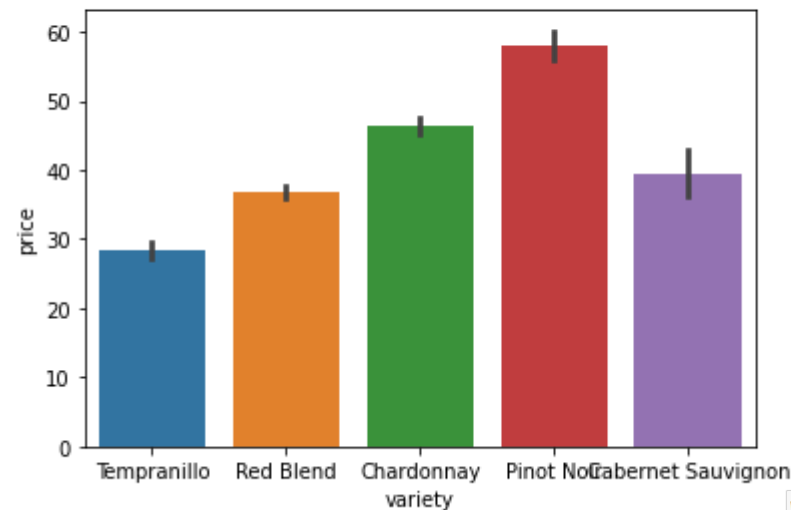
- Use `barplot()` to draw a bar chart grouped by a categorical variable.
- A bar plot shows the **mean** (or other estimator) value.

	country	description	designation	points	price	province	region_1	region_2	variety	winery
51	France	This structured, complex Chardonnay is packed ...	NaN	90	68.0	Burgundy	Chassagne-Montrachet	NaN	Chardonnay	Chartron et Trébouchet
53	France	With its light color and cool feel, this well...	L'Inédit	90	28.0	Loire Valley	Coteaux du Gennois	NaN	Pinot Noir	Clement et Florian Berthier
63	France	L'Homme Mort is a northern extension of the Fo...	L'Homme Mort Premier Cru	91	45.0	Burgundy	Chablis	NaN	Chardonnay	Domaine Chenevrières
66	France	The steely character of a young Chablis is ver...	Fourchaume Premier Cru	91	38.0	Burgundy	Chablis	NaN	Chardonnay	Louis Max
76	France	This wine is bone-dry, although with some age ...	Le Nombre d'Or Brut Nature	91	85.0	Champagne	Champagne	NaN	Chardonnay	Aubry
...	...	...	...	...	...	...	...	...	...	...
150908	France	Another premier cru from Michel Gros, this one...	Aux Brulees	90	65.0	Burgundy	Vosne-Romanée	NaN	Pinot Noir	Michel Gros
150909	France	This is a lovely, fragrant Burgundy, with a sm...	Clos dea Argillieres	89	52.0	Burgundy	Nuits-St-Georges	NaN	Pinot Noir	Daniel Rion
150910	France	Scents of graham cracker and malted milk choco...	NaN	89	38.0	Burgundy	Chambolle-Musigny	NaN	Pinot Noir	Michel Gros
150911	France	This needs a good bit of breathing time, then ...	Les Chaliots	87	37.0	Burgundy	Nuits-St-Georges	NaN	Pinot Noir	Michel Gros
150912	France	The nose is dominated by the attractive scents...	Les Charmes	87	65.0	Burgundy	Chambolle-Musigny	NaN	Pinot Noir	Daniel Rion

4580 rows x 10 columns

```
sns.barplot(data = wine_df, x = "variety", y = "price")
```

```
<AxesSubplot:xlabel='variety', ylabel='price'>
```



The bar height is the average price of each grape variety.

```
wine_df.groupby("variety").mean().price
```

```
variety
Cabernet Sauvignon    39.490964
Chardonnay             46.447540
Pinot Noir            58.053908
Red Blend              36.891409
Tempranillo           28.294762
Name: price, dtype: float64
```

# Heatmap

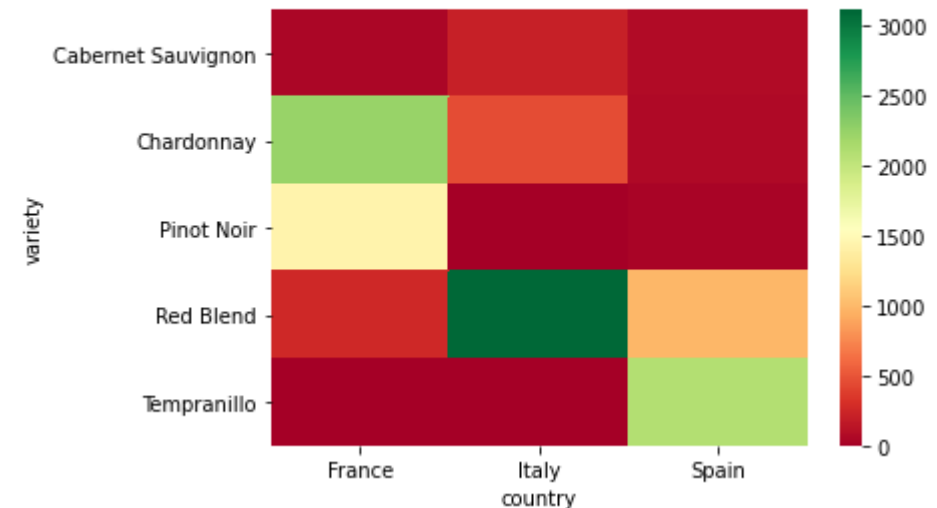
- A heatmap is a two-dimensional visual representation of data using colors, where the colors all represent different values.

```
# data preparation
heatmap_data = pd.crosstab(wine_df["variety"], wine_df["country"])
heatmap_data
```

	country	France	Italy	Spain
variety				
Cabernet Sauvignon		44	214	74
Chardonnay		2247	450	67
Pinot Noir		1446	7	31
Red Blend		267	3111	987
Tempranillo		0	1	2099

Use `crosstab()` to get a cross table of two categorical variables.

```
sns.heatmap(data = heatmap_data, cmap = "RdYlGn")
<AxesSubplot:xlabel='country', ylabel='variety'>
```



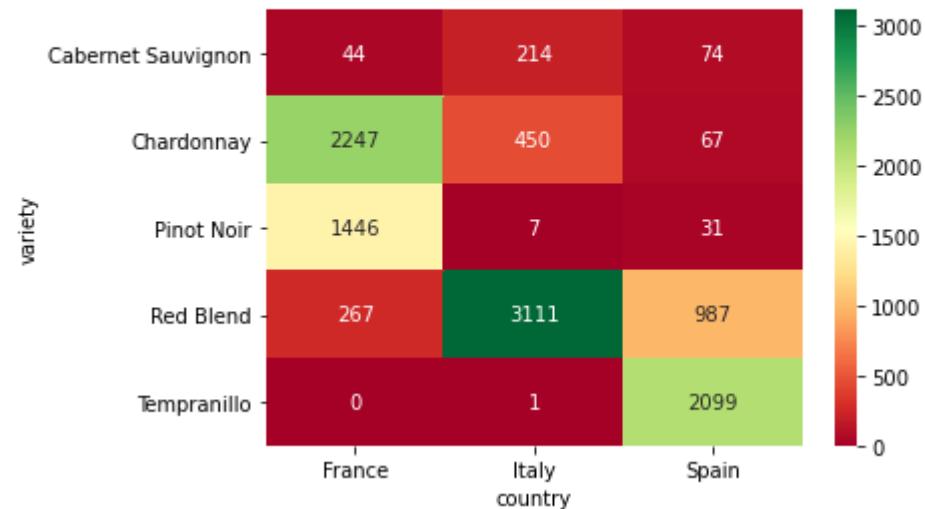
# Heatmap

- Custom style
  - `annot`: If True, show the data value in each cell.
  - `fmt`: String formatting code to use when adding annotations.

```
sns.heatmap(data = heatmap_data, annot = True, fmt = "d", cmap = "RdYlGn")
```

```
<AxesSubplot:xlabel='country', ylabel='variety'>
```

	country	France	Italy	Spain
variety				
Cabernet Sauvignon		44	214	74
Chardonnay		2247	450	67
Pinot Noir		1446	7	31
Red Blend		267	3111	987
Tempranillo		0	1	2099



# Exercise

## Exercise.B

```
titanic_df = sns.load_dataset("titanic", dtype = {"survived": object, "pclass":object})
```

(B.1) Show the first 10 rows of the dataset.

(B.2) Use a count plot to display the number of male and female passengers.

(B.3) Use a count plot to display the number of male and female passengers, and use the column `survived` to divide the data into two subgroups.

Hint: `hue`

(B.4) Create a cross table based on the columns `pclass` and `survived`.

Hint: `pd.crosstab()`

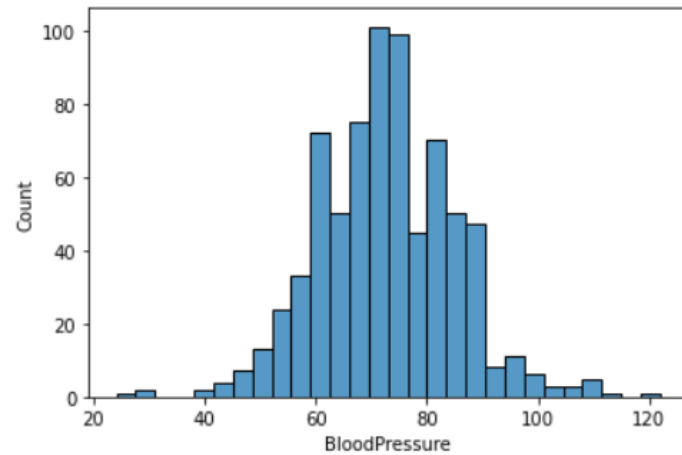
(B.5) Use the result obtained in (B.4) to draw a heatmap.

Setting: `cmap = "coolwarm"`

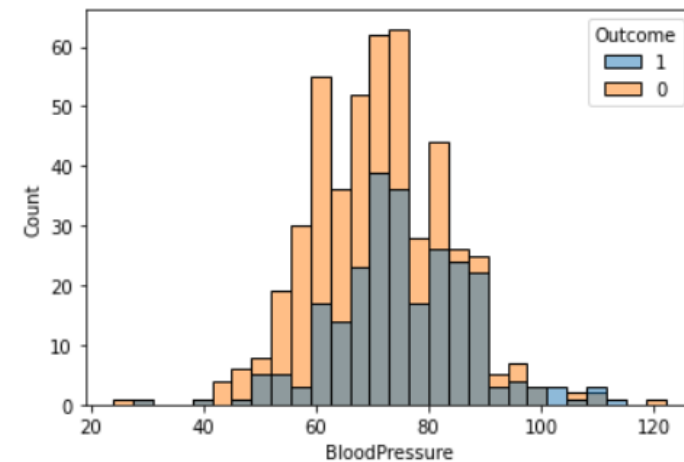
# Histogram

- Use `histplot()` to get the histogram.
- Use `hue` to show the distribution of each category.

```
sns.histplot(data = diabetes_df, x = "BloodPressure")  
<AxesSubplot:xlabel='BloodPressure', ylabel='Count'>
```



```
sns.histplot(data = diabetes_df, x = "BloodPressure", hue = "Outcome")  
<AxesSubplot:xlabel='BloodPressure', ylabel='Count'>
```



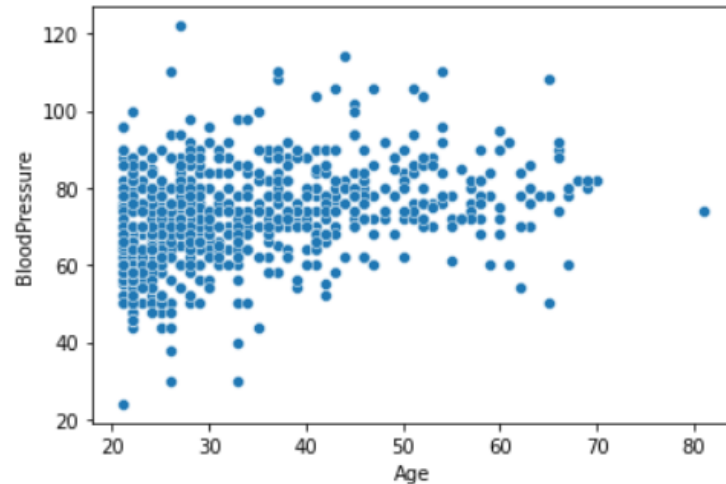


# Scatter plot

- Use `scatterplot()` to show the relationship between two variables.
- Use `hue` to show the scatter plot for each category.

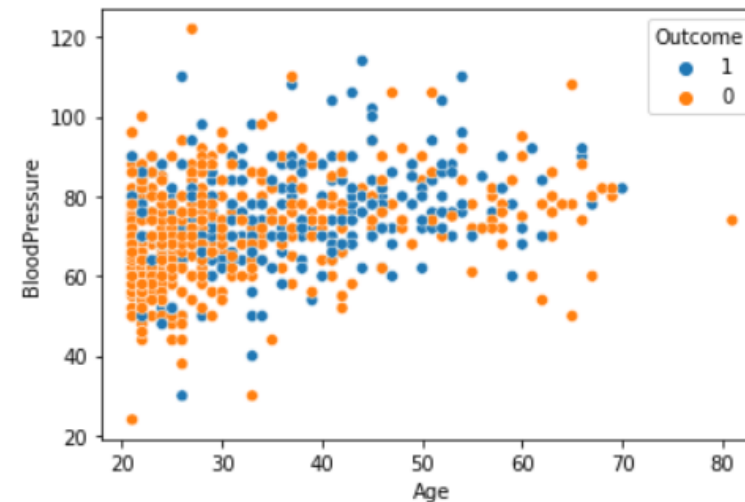
```
sns.scatterplot(data = diabetes_df, x = "Age", y = "BloodPressure")
```

```
<AxesSubplot:xlabel='Age', ylabel='BloodPressure'>
```



```
sns.scatterplot(data = diabetes_df, x = "Age", y = "BloodPressure",  
               hue = "Outcome")
```

```
<AxesSubplot:xlabel='Age', ylabel='BloodPressure'>
```

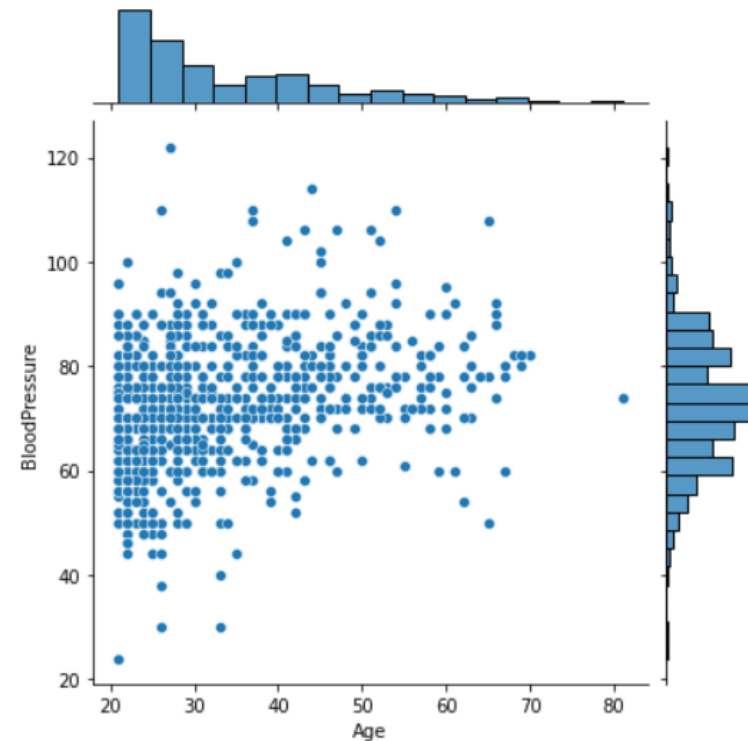


# Jointplot

- Use `jointplot()` to show the individual distribution and the relationship between two variables.

```
sns.jointplot(data = diabetes_df, x = "Age", y = "BloodPressure")
```

```
<seaborn.axisgrid.JointGrid at 0x19793e56640>
```

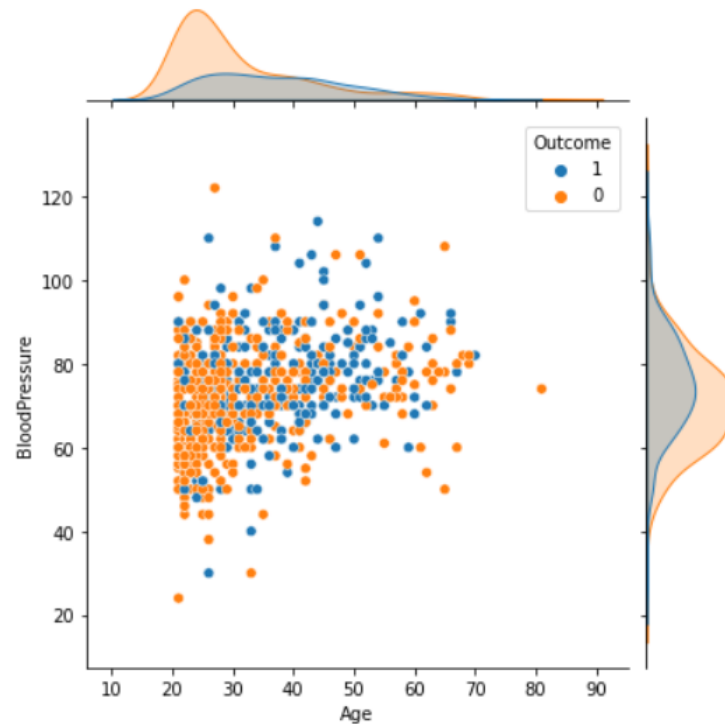


# Jointplot

- Use the argument **hue** to specify the categorical variable used for color encoding.

```
sns.jointplot(data = diabetes_df, x = "Age", y = "BloodPressure", hue = "Outcome")
```

```
<seaborn.axisgrid.JointGrid at 0x19793df2e20>
```

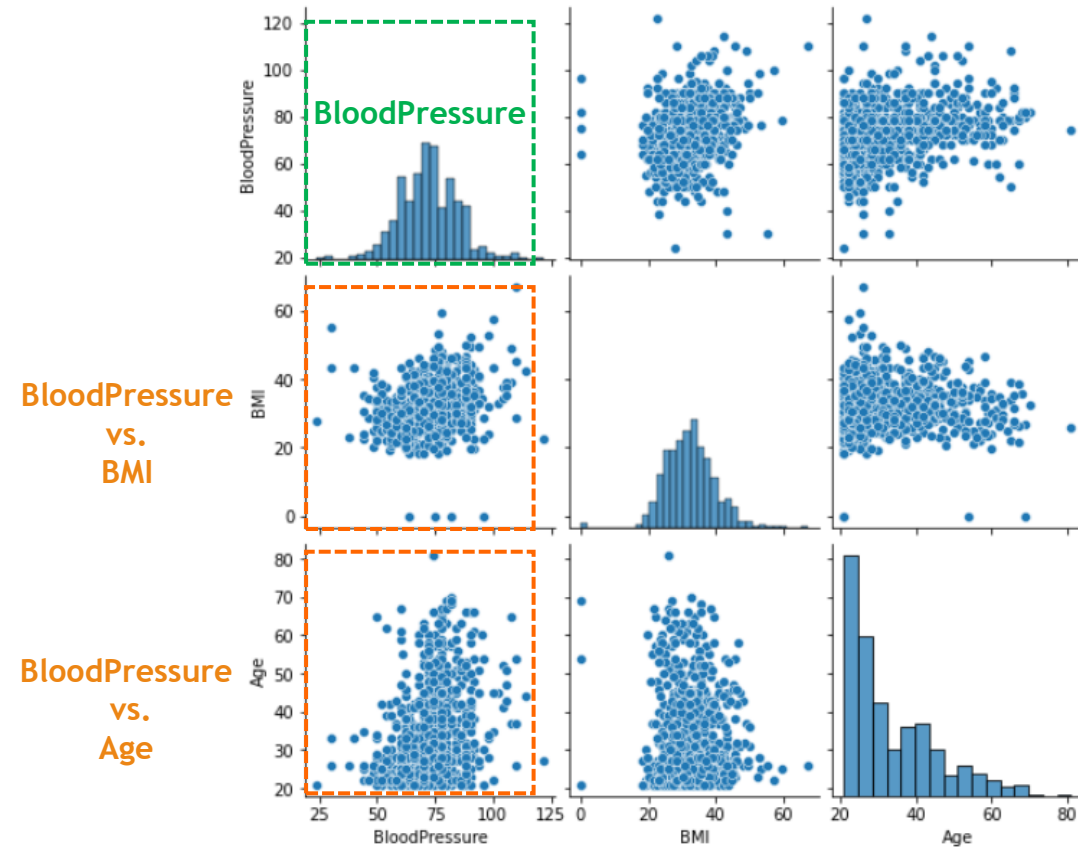


# Pair plot

- Use `Pairplot()` to plot pairwise relationships in a dataset.
- The **diagonal plots** are univariate distribution plot (histogram) of the data in each column.

```
sns.pairplot(data = diabetes_df.loc[:,["BloodPressure","BMI", "Age"]])
```

<seaborn.axisgrid.PairGrid at 0x1de936fda30>



# Exercise

---

## Exercise.C

```
diamond_df = sns.load_dataset("diamonds")
```

(C.1) Show the first 5 rows of the dataset.

(C.2) Show the distribution of the price.

Hint: `histplot()`

(C.3) Use a scatter plot to show the relationship between `carat` and `price` of the best clarity diamonds.

Hint: Select a subset by using `diamond_df.clarity == "IF"`

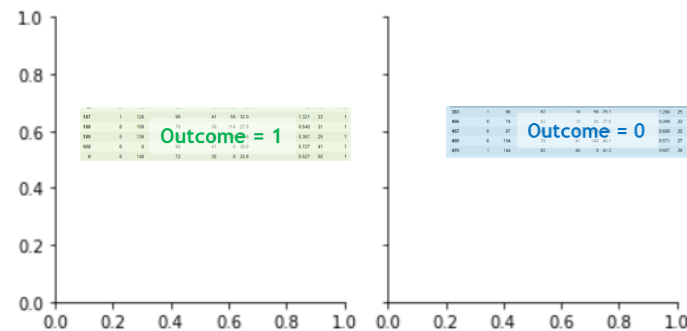
(C.4) Use a join plot to show the relationship between `carat` and `price` of the best clarity diamonds.

# FacetGrid

- FacetGrid is basically a grid of subplots.
- The method `FacetGrid()` return a `FacetGrid` object which stores some information on how you want to break down your data visualization.

## Step1: Initialize the grid

```
fig = sns.FacetGrid(data = diabetes_df, col = "Outcome")
```



```
type(fig)
```

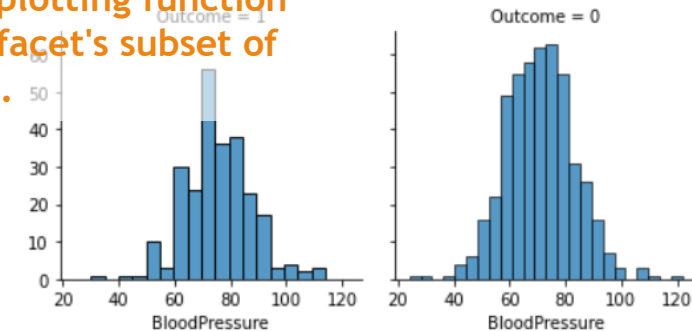
```
seaborn.axisgrid.FacetGrid
```

## Step2: Draw a plot on every facet

```
fig = sns.FacetGrid(data =diabetes_df, col="Outcome")  
fig.map(sns.histplot, "BloodPressure")
```

```
<seaborn.axisgrid.FacetGrid at 0x25650e88730>
```

Apply a plotting function to each facet's subset of the data.

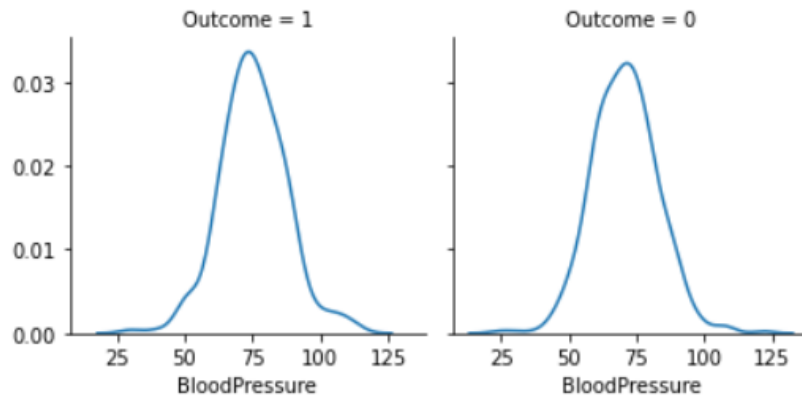


# FacetGrid

- Change chart type

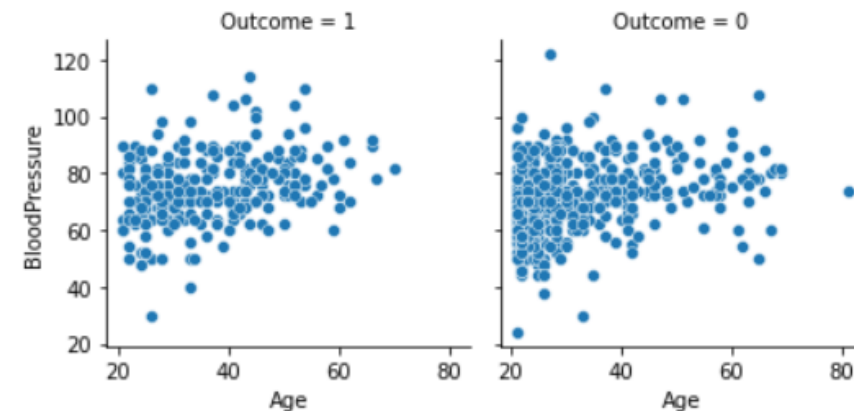
```
fig = sns.FacetGrid(data = diabetes_df, col="Outcome")  
fig.map(sns.kdeplot, "BloodPressure")
```

<seaborn.axisgrid.FacetGrid at 0x25650c48e80>



```
fig = sns.FacetGrid(data = diabetes_df, col = "Outcome")  
fig.map(sns.scatterplot, "Age", "BloodPressure")
```

<seaborn.axisgrid.FacetGrid at 0x25650d6ac10>



# Exercise

## Exercise.D

(D.1) Use the dataframe `diamond_df` in (C.1). Draw a price histogram for each `cut` category.

Hint: `col="cut"`

