



Pandas 4 - Data Visualization



Data analysis process

Data Understanding

- Descriptive statistics
- Types of data (numerical/categorical)

Data Preprocessing

- Subset selection
- Data consolidation
- Missing data handling

Calculation (Modeling)

- Derived variables
- Data aggregation

Data Visualization

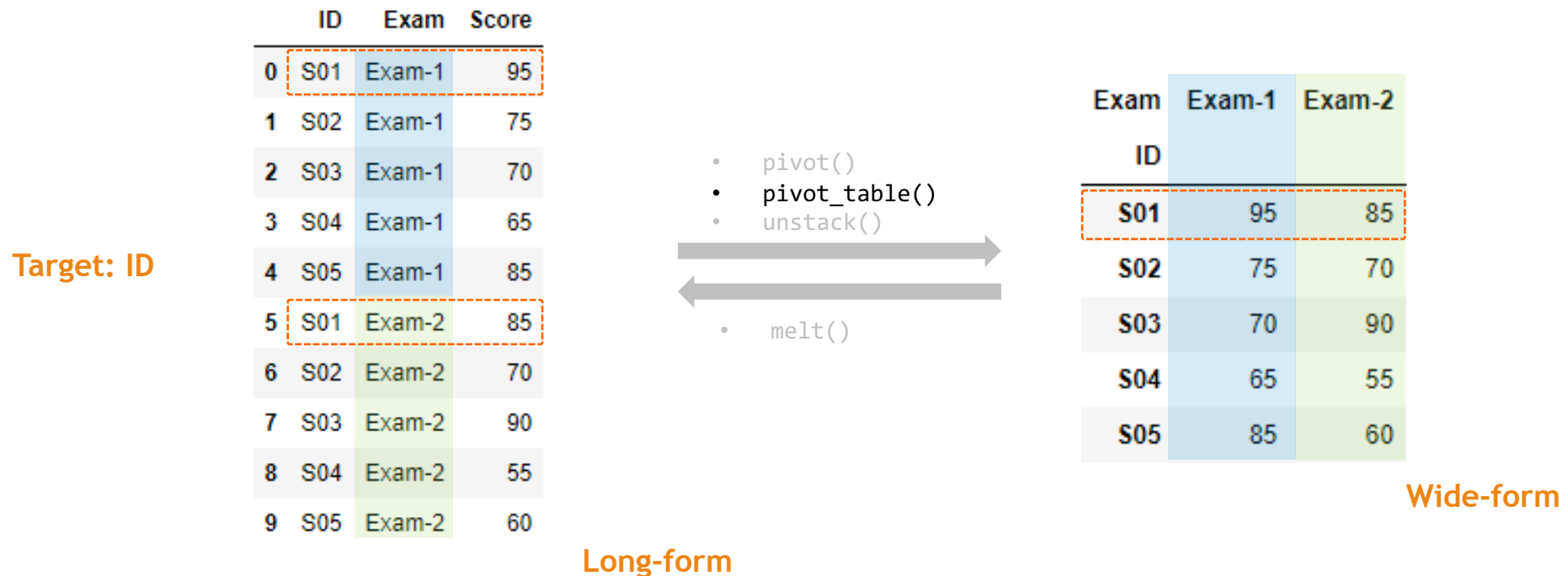
- Univariate chart
- Bivariate chart
- Multivariate chart

Outline

- Reshape DataFrame for visualization
- X-axis with categorical data
 - (Line chart)
 - Bar chart
 - Area chart
 - Pie chart
- Numerical data
 - Histogram
 - Scatter plot
 - Hexagon plot

Wide-form and Long-form

- **Long-form:** Each row is one time point per target. The data of a target can have **multiple rows**.
- **Wide-form:** A target's repeated responses will be in a **single row**, and each response is in a separate column.



Pivot_table

- Use `pivot_table()` to reshaped a DataFrame by passing arguments: index, columns and values. (By default, aggfunc = mean.)

	ID	Exam	Score
0	S01	Exam-1	95
1	S02	Exam-1	75
2	S03	Exam-1	70
3	S04	Exam-1	65
4	S05	Exam-1	85
5	S01	Exam-2	85
6	S02	Exam-2	70
7	S03	Exam-2	90
8	S04	Exam-2	55
9	S05	Exam-2	60

Long-form

```
score_df.pivot_table(index = "ID", columns = "Exam", values = "Score")
```

Exam	Exam-1	Exam-2
ID		
S01	95	85
S02	75	70
S03	70	90
S04	65	55
S05	85	60

Wide-form

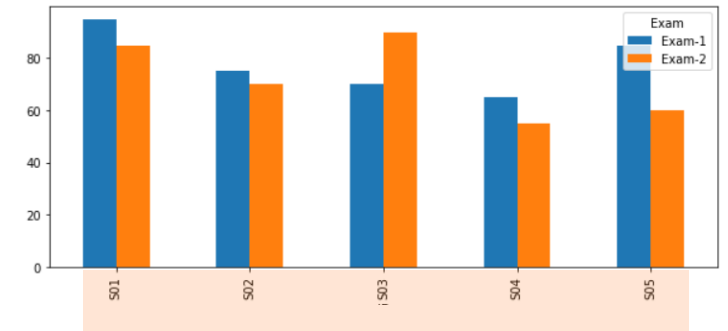
Pivot_table

- Change the target.

```
score_df.pivot_table(index = "ID", columns = "Exam", values = "Score")
```

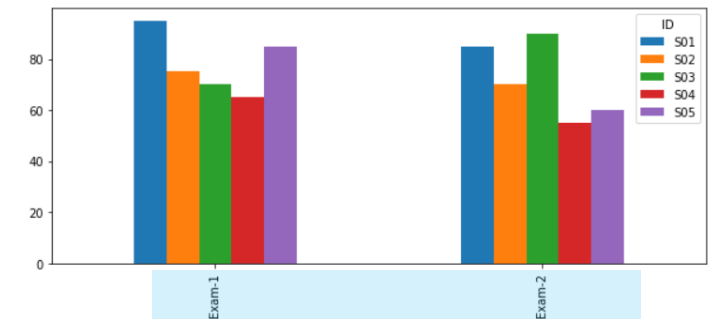
	ID	Exam	Score
0	S01	Exam-1	95
1	S02	Exam-1	75
2	S03	Exam-1	70
3	S04	Exam-1	65
4	S05	Exam-1	85
5	S01	Exam-2	85
6	S02	Exam-2	70
7	S03	Exam-2	90
8	S04	Exam-2	55
9	S05	Exam-2	60

Exam	Exam-1	Exam-2
ID		
S01	95	85
S02	75	70
S03	70	90
S04	65	55
S05	85	60



```
score_df.pivot_table(index = "Exam", columns = "ID", values = "Score")
```

ID	S01	S02	S03	S04	S05
Exam					
Exam-1	95	75	70	65	85
Exam-2	85	70	90	55	60



Unstack



- Use `unstack()` to reshape a dataframe/series derived from a groupby object.
- By default, `unstack()` converts the inner-most row level to column level.

	Product	Quarter	Month	Sales
0	A	Q1	Jan	67
1	A	Q1	Feb	57
2	A	Q1	Mar	87
3	A	Q2	Apr	50
4	A	Q2	May	97
5	A	Q2	Jun	68
6	B	Q1	Jan	78
7	B	Q1	Feb	102
8	B	Q1	Mar	113
9	B	Q2	Apr	98
10	B	Q2	May	80
11	B	Q2	Jun	84



```
sales_df.groupby(["Product", "Quarter"]).Sales.sum()
```

Product	Quarter	
A	Q1	211
	Q2	215
B	Q1	293
	Q2	262

Name: Sales, dtype: int64

Long-form

```
sales_df.groupby(["Product", "Quarter"]).Sales.sum().unstack()
```

Quarter	Q1	Q2
Product		
A	211	215
B	293	262

wide-form

Melt



- Use `melt()` to convert a DataFrame from wide-form to long-form.

Step1: Convert "ID" from index to column

Exam	Exam-1	Exam-2
ID		
S01	95	85
S02	75	70
S03	70	90
S04	65	55
S05	85	60

wide-form

```
wide_df.reset_index(inplace = True)
wide_df
```

Exam	ID	Exam-1	Exam-2
0	S01	95	85
1	S02	75	70
2	S03	70	90
3	S04	65	55
4	S05	85	60

Step2: Convert wide-form to long-form

```
long_df = pd.melt(wide_df,
                   id_vars = "ID",
                   var_name="Exam",
                   value_name="Score")
long_df
```

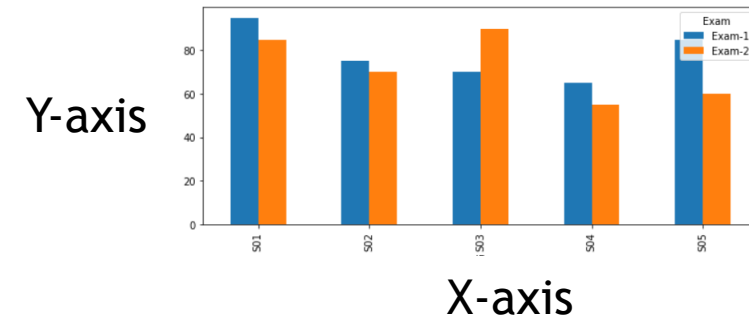
	ID	Exam	Score
0	S01	Exam-1	95
1	S02	Exam-1	75
2	S03	Exam-1	70
3	S04	Exam-1	65
4	S05	Exam-1	85
5	S01	Exam-2	85
6	S02	Exam-2	70
7	S03	Exam-2	90
8	S04	Exam-2	55
9	S05	Exam-2	60

Long-form

- Melt: <https://pandas.pydata.org/docs/reference/api/pandas.melt.html>

Outline

- Reshape DataFrame for visualization
- X-axis with categorical data
 - (Line chart)
 - Bar chart
 - Area chart
 - Pie chart
- Numerical data
 - Histogram
 - Scatter plot
 - Hexagon plot



Packages - packages 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

Line chart - Series

- Both Series and DataFrame have a `plot()` method to make some basic plot types. By default, `plot()` makes line charts.
- A line chart is usually used to visualize the trend of data over a period of time.

```
series_A = pd.Series([67, 57, 87, 50, 97, 68],  
                     index = ["Jan", "Feb", "Mar", "Apr", "May", "Jun"])
```

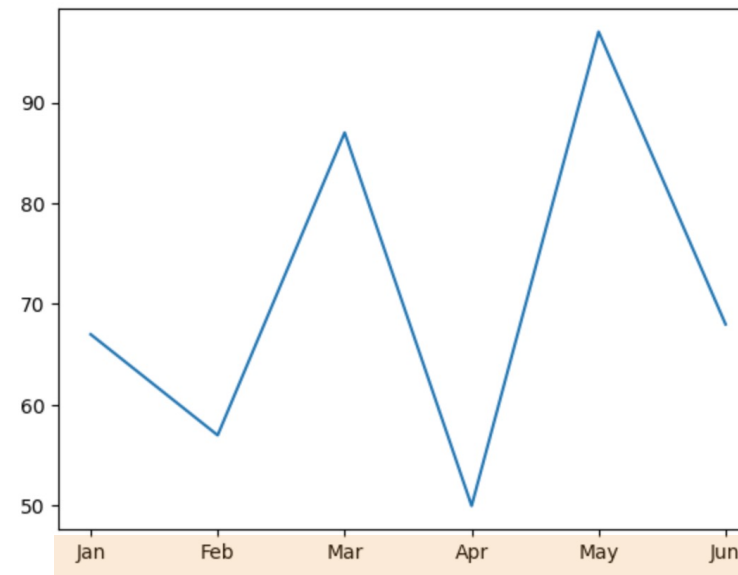
series_A

Jan	67
Feb	57
Mar	87
Apr	50
May	97
Jun	68

dtype: int64

```
series_A.plot()
```

<Axes: >



Use index as ticks
for x-axis.

Line chart - DataFrame

- Use “x” and “y” to specify the columns used for plotting.

	month	sales_A	sales_B
0	Jan	67	78
1	Feb	57	102
2	Mar	87	113
3	Apr	50	98
4	May	97	80
5	Jun	68	84

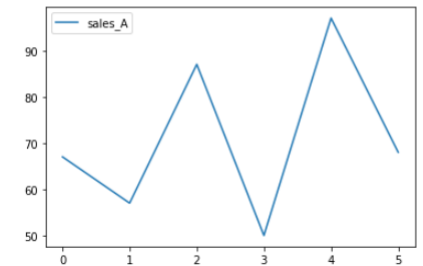
```
product_df.plot(x = "month", y = "sales_A")
```

<AxesSubplot:xlabel='month'>



```
product_df.plot(y = "sales_A")
```

<AxesSubplot:>



If “x” is not specified, the index of the DataFrame is used.

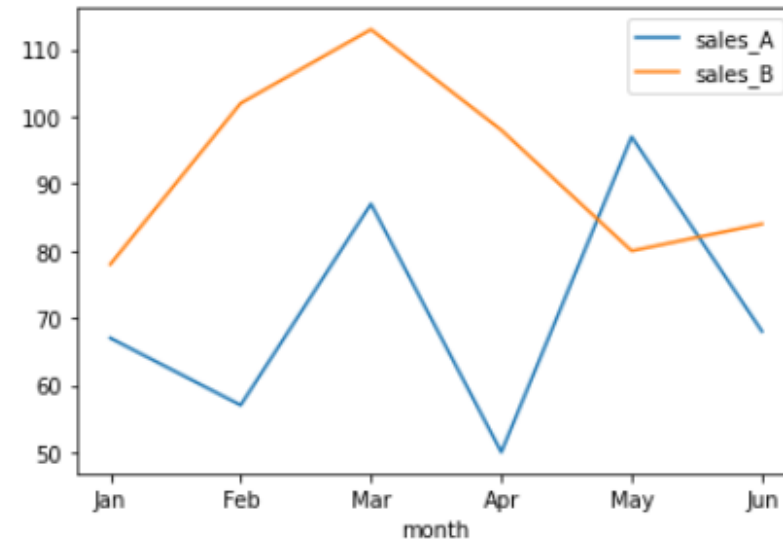
Line chart - Multiple lines

- Pass a list of column names to the argument “y” to plot multiple lines.

	month	sales_A	sales_B
0	Jan	67	78
1	Feb	57	102
2	Mar	87	113
3	Apr	50	98
4	May	97	80
5	Jun	68	84

```
product_df.plot(x = "month", y= ["sales_A","sales_B"])
```

<AxesSubplot:xlabel='month'>



➔ Product B's sales exceeded Product A's in all months except May.

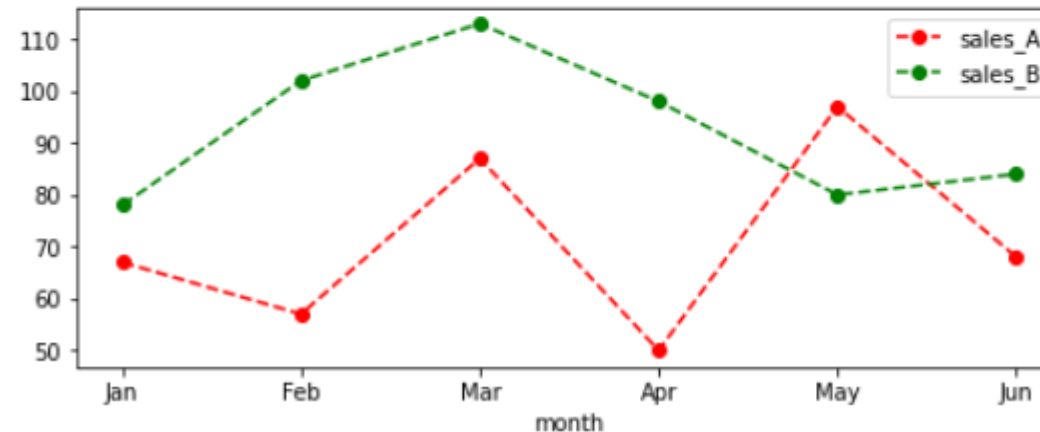
Line chart - Custom style

- Use some arguments to change the style.

	month	sales_A	sales_B
0	Jan	67	78
1	Feb	57	102
2	Mar	87	113
3	Apr	50	98
4	May	97	80
5	Jun	68	84

```
product_df.plot(x = "month",  
                y = ["sales_A", "sales_B"],  
                marker = "o",  
                color = ["red", "green"],  
                linestyle = 'dashed',  
                figsize = (8,3))
```

<AxesSubplot:xlabel='month'>



- Marker: https://matplotlib.org/stable/api/markers_api.html#module-matplotlib.markers
- Color: https://matplotlib.org/stable/gallery/color/named_colors.html
- Linestyle: https://matplotlib.org/stable/gallery/lines_bars_and_markers/linestyles.html
- Others: https://matplotlib.org/stable/api/_as_gen/matplotlib.pyplot.plot.html

Exercise

Exercise.A

(A.1) Given the dataframe `expense_df` . Convert the dataframe to the following format (wide-form) and store the result in a new variable named `expense_df_wide` .

(A.2) Use the dataframe `expense_df_wide` obtained in (A.1). Draw a multiple line chart to show the monthly groceries and transportation expenses.

(A.3) Import dataset `fashion.csv` . Show the first five rows.

(A.4) Show the sales trends of `Tiger_of_Sweden` with a line chart.

(A.5) Show the sales trends of `Eton` , `Levi_s` , and `Tiger_of_Sweden` with a multiple line chart.

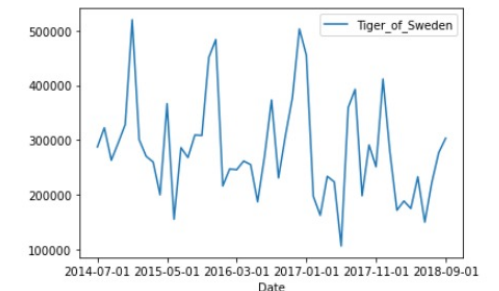
Settings: Use `marker = "D"` , `figsize = (12,4)` , `title = "Monthly Sales"` , `ylabel = "Sales"` .

(A.1)

	month	expense	category
0	01-2022	3050	grocery
1	02-2022	2800	grocery
2	03-2022	2750	grocery
3	04-2022	2300	grocery
4	05-2022	3150	grocery
5	06-2022	2900	grocery
6	01-2022	1050	transportation
7	02-2022	900	transportation
8	03-2022	1150	transportation
9	04-2022	1850	transportation
10	05-2022	1250	transportation
11	06-2022	950	transportation

category	grocery	transportation
month		
01-2022	3050	1050
02-2022	2800	900
03-2022	2750	1150
04-2022	2300	1850
05-2022	3150	1250
06-2022	2900	950

(A.4)



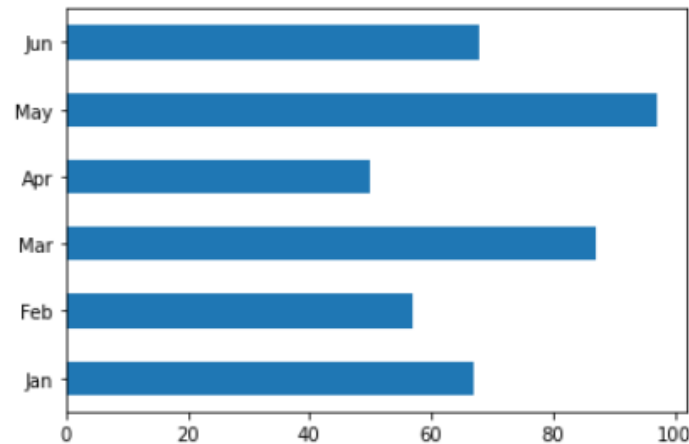
Bar chart - Series

- A bar chart is used to compare values of different categories.
 - Use `kind = "bar"` to plot vertical bar chart.
 - Use `kind = "barh"` to plot horizontal bar chart.

```
Jan    67  
Feb    57  
Mar    87  
Apr    50  
May    97  
Jun    68  
dtype: int64
```

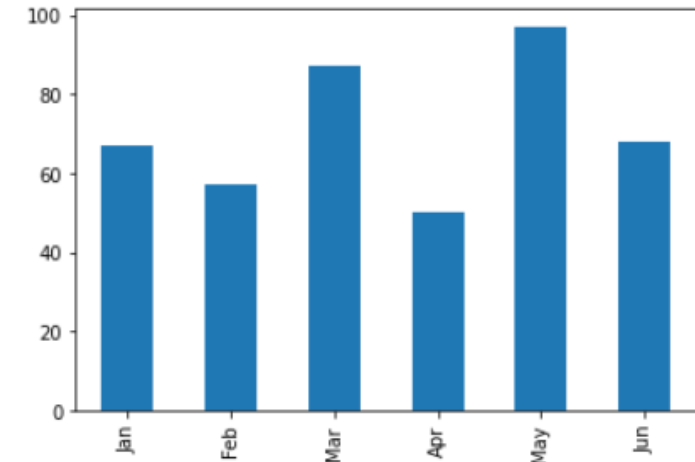
```
series_A.plot(kind = 'barh')
```

<AxesSubplot:>



```
series_A.plot(kind = 'bar')
```

<AxesSubplot:>



Bar chart - DataFrame

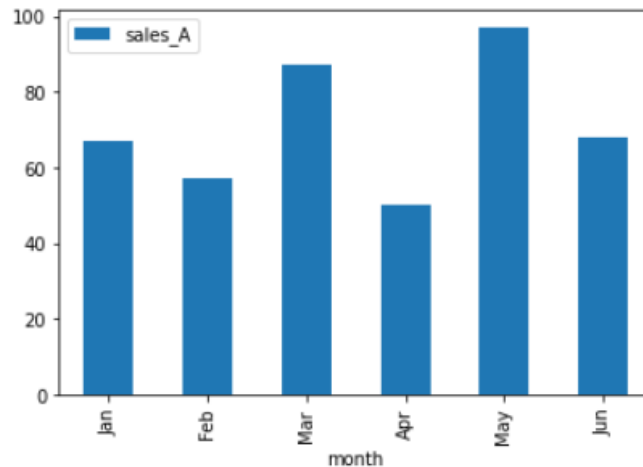
- Use “x” and “y” to specify the columns used for plotting.

Pass a list of column names to the argument “y” to plot multiple bars.

```
product_df.plot(kind = "bar", x = "month", y = "sales_A")
```

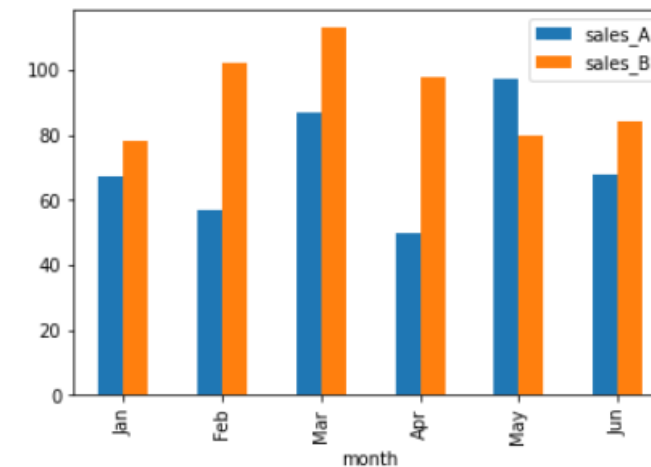
<AxesSubplot:xlabel='month'>

	month	sales_A	sales_B
0	Jan	67	78
1	Feb	57	102
2	Mar	87	113
3	Apr	50	98
4	May	97	80
5	Jun	68	84



```
product_df.plot(kind = "bar", x = "month", y = ["sales_A", "sales_B"])
```

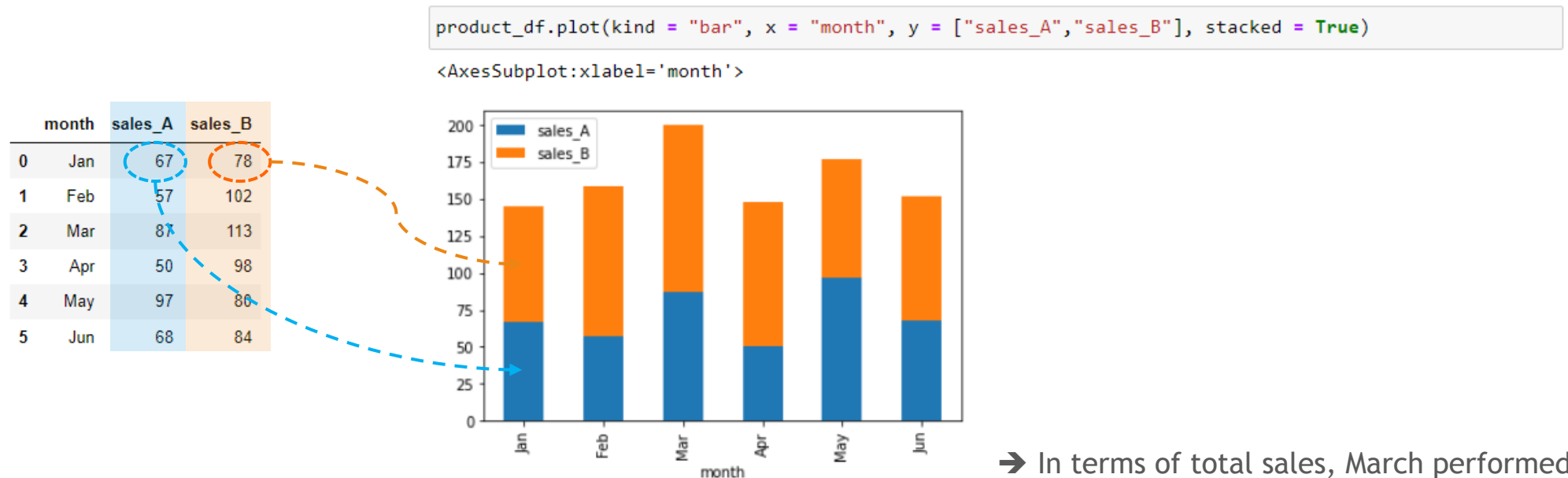
<AxesSubplot:xlabel='month'>



Bar chart - Stacked bar chart

- Stacked bar chart

- Each bar is stacked by multiple data series.
- Stacked bar charts can be used to break down and compare parts of the whole.



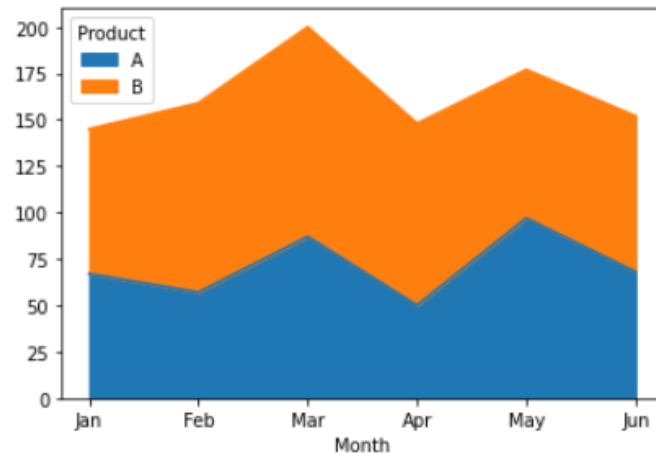
➔ In terms of total sales, March performed best, followed by May.

Area chart

- Area charts are similar to line charts, except that the area below the line is filled with color, making it easier to understand the cumulative value.
- Use `kind = "area"` to plot an area chart. By default, `stacked = True`.

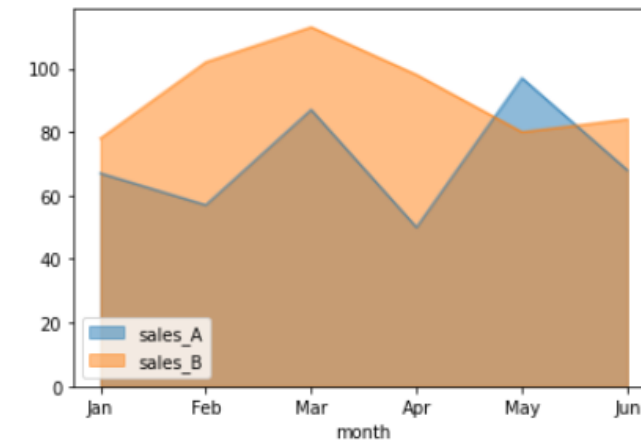
```
product_df.plot(kind = "area",  
                x = "month",  
                y = ["sales_A", "sales_B"])
```

<AxesSubplot:xlabel='Month'>



```
product_df.plot(kind = "area",  
                x = "month",  
                y = ["sales_A", "sales_B"],  
                stacked = False)
```

<AxesSubplot:xlabel='month'>

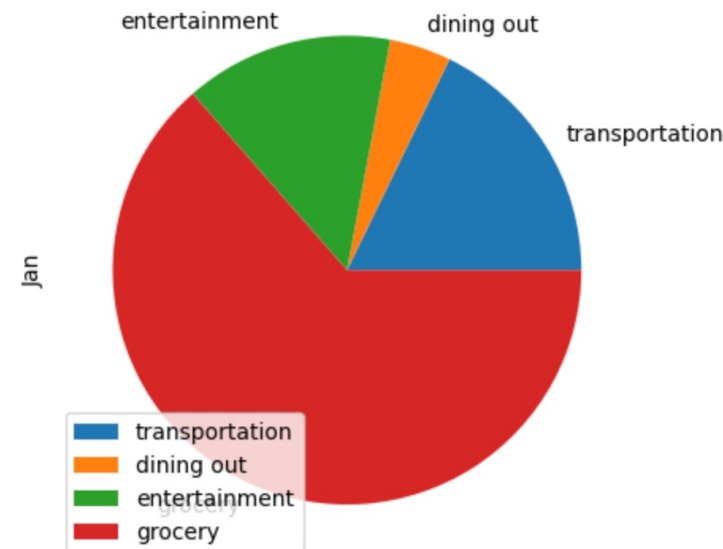


Pie chart

- A pie chart is used to show the proportion of each category to the whole.

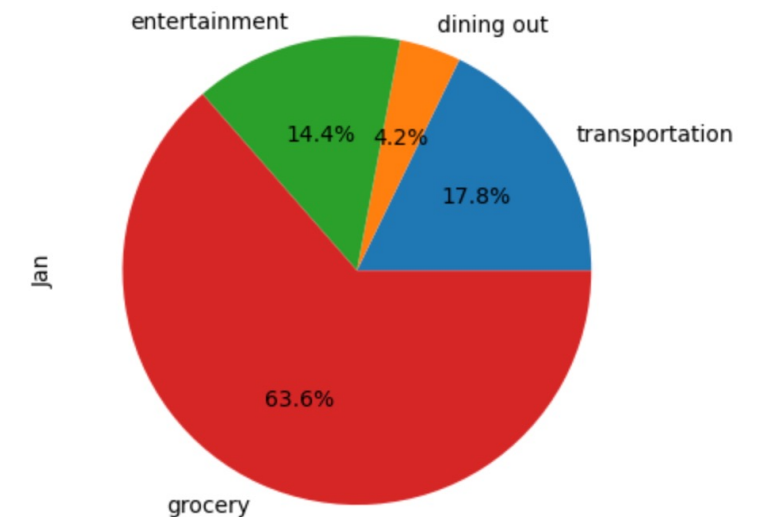
	Jan	Feb	Mar
transportation	1050	1750	1150
dining out	250	850	450
entertainment	850	1050	950
grocery	3750	3050	3250

```
spend_df.plot(kind = "pie", y = "Jan")  
<Axes: ylabel='Jan'>
```



```
spend_df.plot(kind = "pie", y = "Jan",  
              autopct='%.1f%%',  
              legend = False)
```

<Axes: ylabel='Jan'>



BI

Autopct = `%.1f%%`
(Auto percentage)

- `'%'` is part of the formatting command `%.1f` and is used to specify the format of a floating-point number.
- `'%'` is an escape character that indicates that the following `'%'` should be treated as a literal `'%'` symbol.
- `'%%'` is treated as a literal `'%'` symbol.



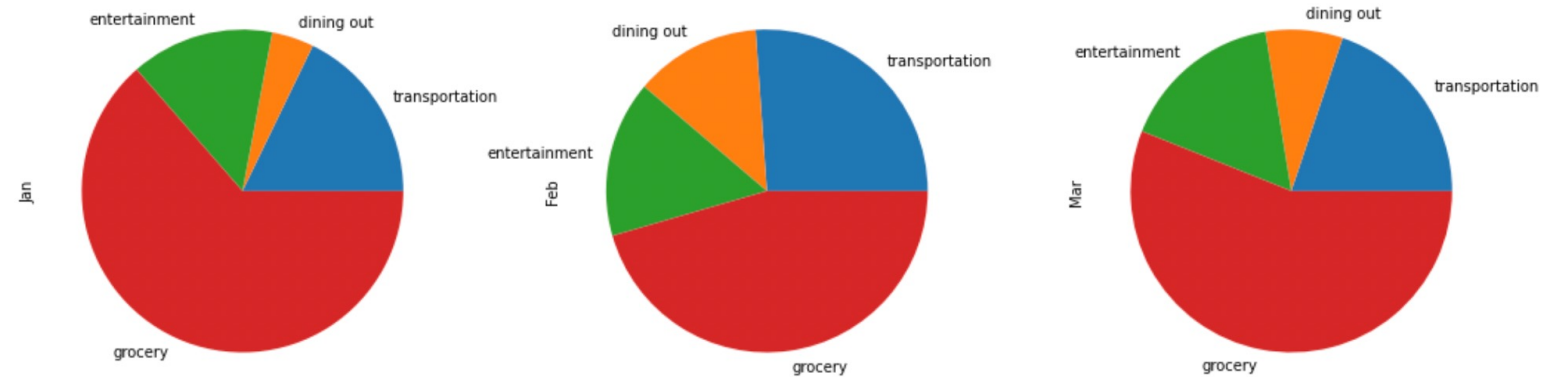
Pie chart



- Use `subplots=True` to plot a pie chart for each numerical column.

	Jan	Feb	Mar
transportation	1050	1750	1150
dining out	250	850	450
entertainment	850	1050	950
grocery	3750	3050	3250

```
spend_df.plot(kind = "pie", subplots = True, figsize=(18,5), legend = False);
```



Exercise

(B.1) Import dataset `parks.csv` . Show the first five rows.

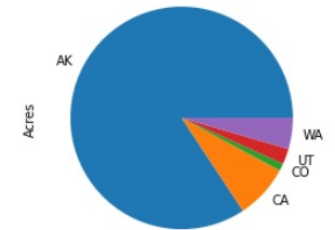
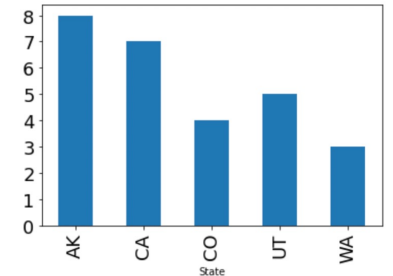
(B.2) Select the national parks in the following five states and keep columns `Park Name` , `State` , and `Acres` . Use this subset to answer the following questions.

State: CA, CO, UT, AK, WA

(B.3) Count the number of national parks in each state. Display the result using a bar graph.

Hint: (1) Group data using column "State". (2) The x-axis shows each state, and each bar is the number of national parks in each state.

(B.4) Calculate the total area of national parks in each state. Display the result using a pie chart.



Outline

- Reshape DataFrame for visualization
- X-axis with categorical data
 - (Line chart)
 - Bar chart
 - Area chart
 - Pie chart
- Numerical data
 - Histogram
 - Scatter plot
 - Hexagon plot

Histogram

- A histogram is used to display the distribution of numerical data.
 - Step1: Divide the entire range of values into a series of intervals.
 - Step2: Count how many values fall into each interval.

	Pregnancies	Glucose	BloodPressure	SkinThickness
0	6	148	72	35
1	1	85	66	29
2	8	183	64	0
3	1	89	66	23
4	0	137	40	35
...
763	10	101	76	48
764	2	122	70	27
765	5	121	72	23
766	1	126	60	0
767	1	93	70	31

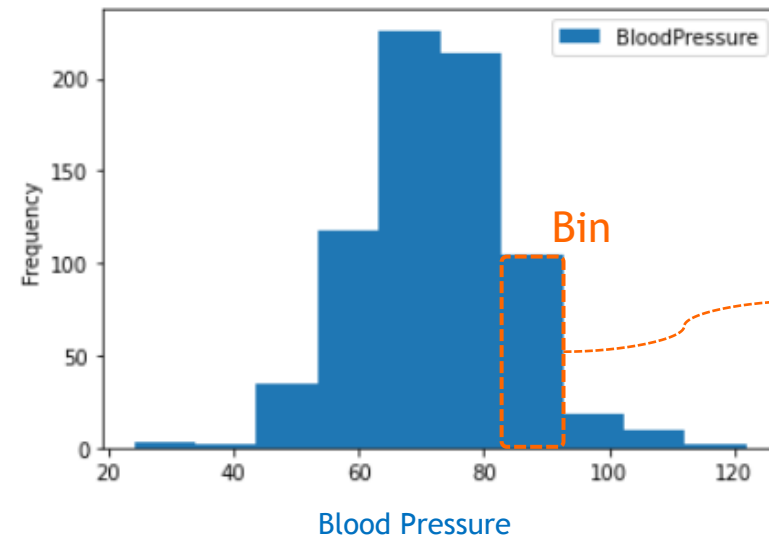
By default, the number of bins is 10.

width of bins
= (max - min)/number of bins
= (122-24)/10 = 9.8

Number of people

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

<AxesSubplot:ylabel='Frequency'>



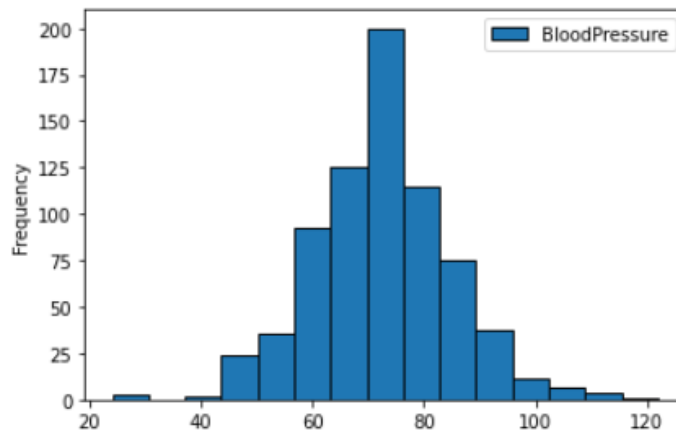
There were 105 people with blood pressure between 82.8 and 92.6.

Histogram - custom bins

- Use the argument “bins” to customize the number of bins.
 - Integer: bins = 15.
 - A sequence of bin edges: bins = [0,5,10,...,130]

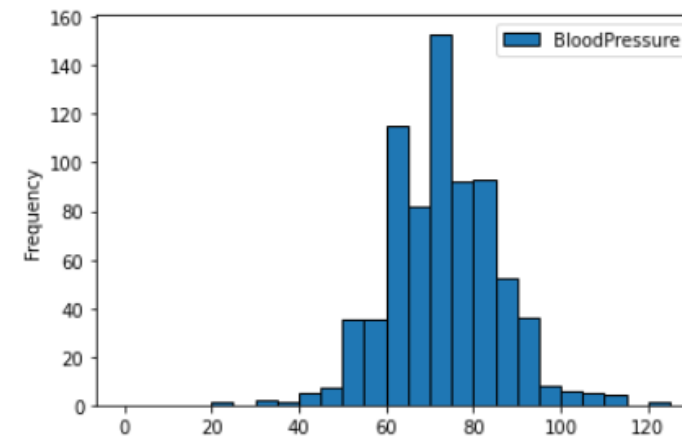
```
diabetes_df.plot(kind = "hist",  
                 y = "BloodPressure",  
                 bins = 15,  
                 edgecolor = "black")
```

<AxesSubplot:ylabel='Frequency'>



```
diabetes_df.plot(kind = "hist",  
                 y = "BloodPressure",  
                 bins = range(0,130,5),  
                 edgecolor = "black")
```

<AxesSubplot:ylabel='Frequency'>

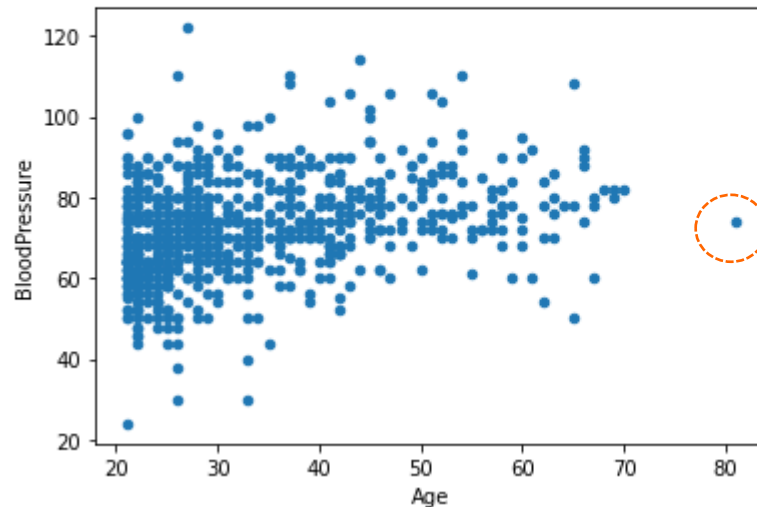


Scatter plot

- Scatter plots are used to observe the relationship between two variables.
 - Each dot indicates an individual data point (observation).

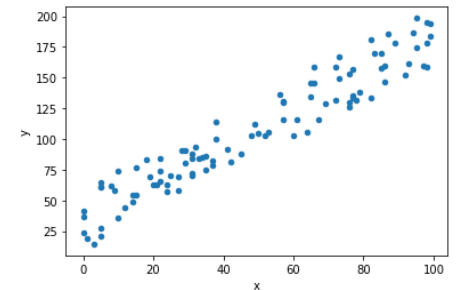
```
diabetes_df.plot(kind = "scatter", x = "Age", y = "BloodPressure")
```

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

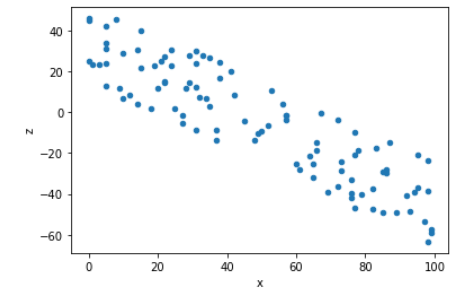


Age = 81,
BloodPressure = 74

Positive correlation



Negative correlation

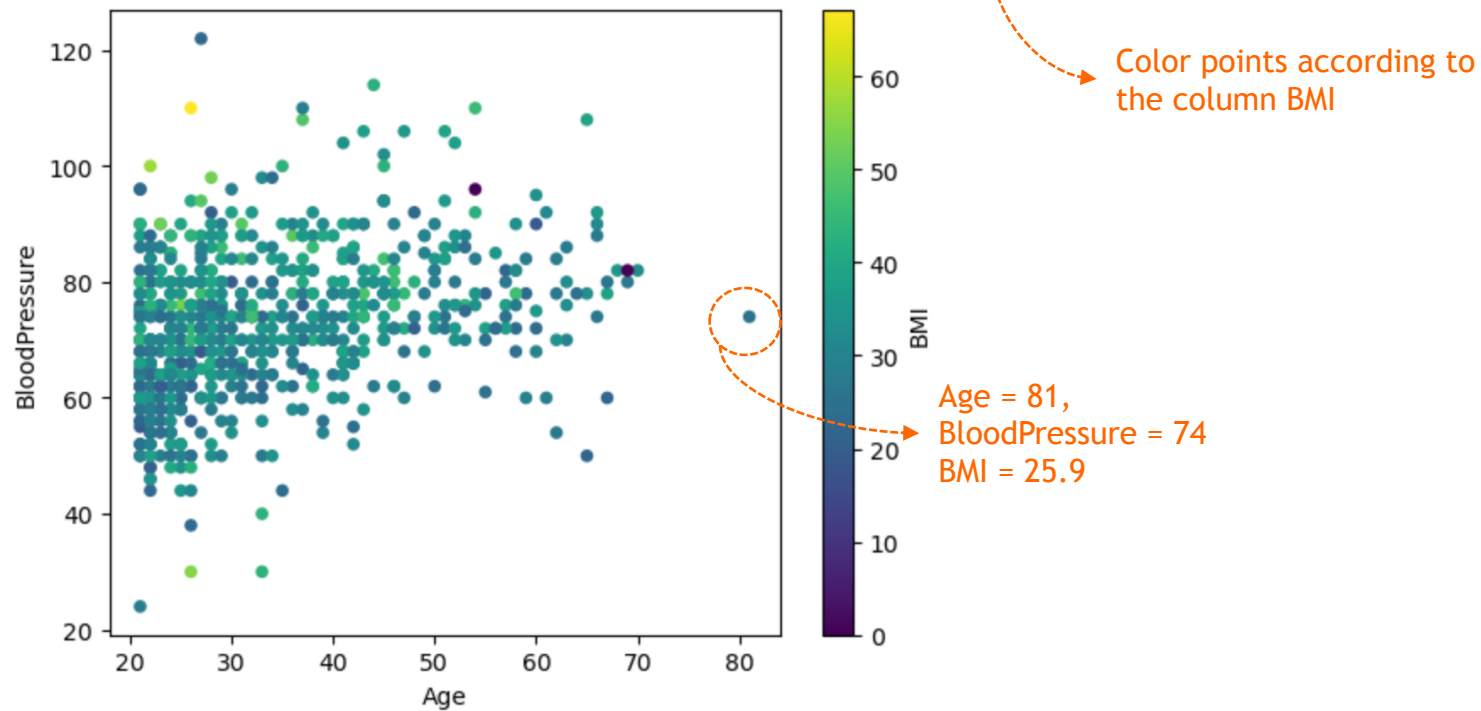


Scatter plot

- Color points based on the third variable.

```
diabetes_df.plot(kind = "scatter", x = "Age", y = "BloodPressure", c = "BMI", cmap = "viridis")
```

<Axes: xlabel='Age', ylabel='BloodPressure'>



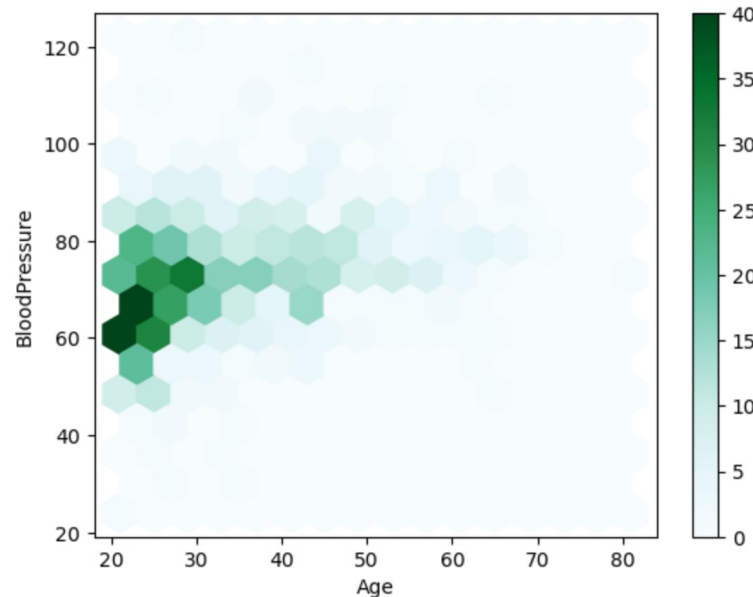
Hexagon plot



- A hexagon plot combines nearby data points into a hexagon, and then displays the **density** (the number of data points) in color.
- Hexagon plots can solve the problem that many points begin to overlap.

```
diabetes_df.plot(kind = "hexbin", x = "Age", y = "BloodPressure", gridsize = 15)
```

<Axes: xlabel='Age', ylabel='BloodPressure'>



Gridsize: The number of hexagons in the x-direction.

Exercise

Exercise.C

(C.1) Import dataset `wine.csv` and set the first column as the index. Display the first 5 rows.

Hint: `index_col = [0]`

(C.2) Select a subset that satisfies the following two conditions. Use this subset for the following tasks.

- Select wines (rows) from Spain, Italy or France (use column `country`).
- Select wines (rows) with a price of less than 200 (use column `price`).

(C.3) Use a histogram to show the price distribution of French wines.

Hint: Use column `price`.

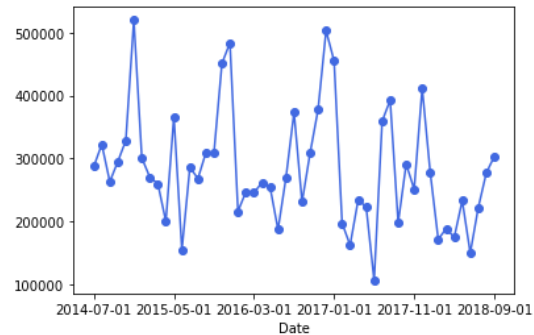
(C.4) Use a scatter plot to show the relationship between price and the points received in the review.

Hint: Use column `price` and `points`.

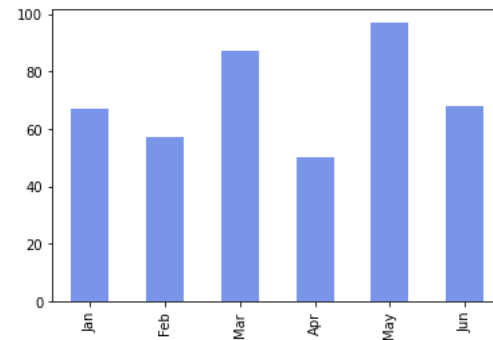
Summary

Univariate chart

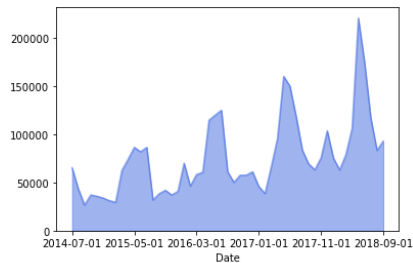
Line chart



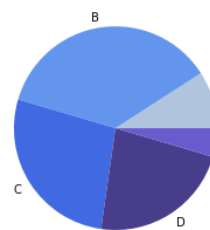
Bar chart



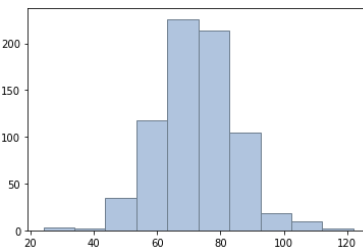
Area chart



Pie chart

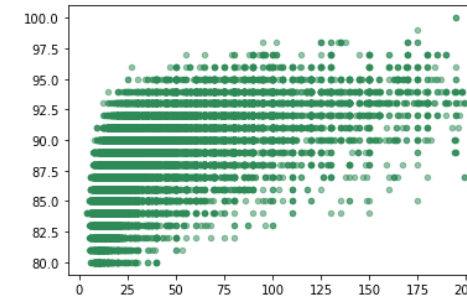


Histogram

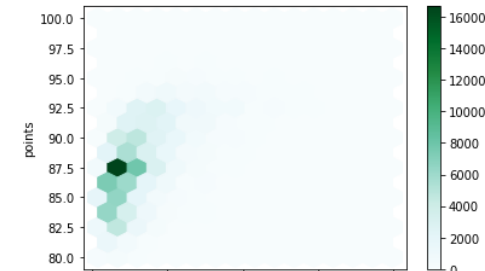


Bivariate chart

Scatter plot

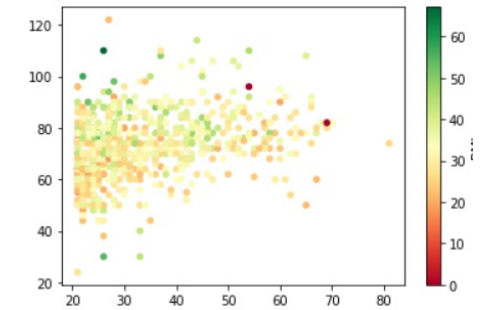


Hexagon plot



Multivariate chart

Scatter plot



Heatmap

