



Exploratory Data Analysis (P1)

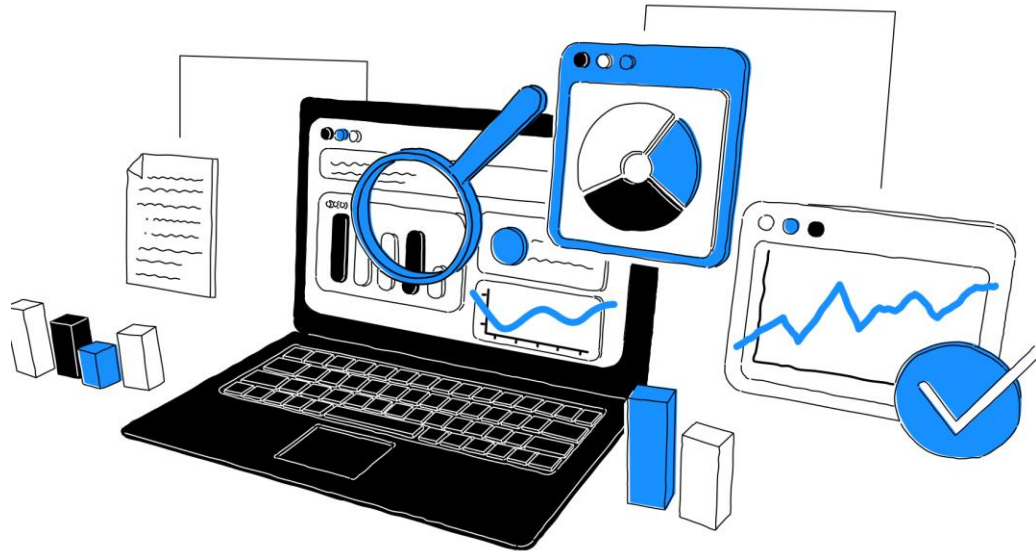
Nguyen Ngoc Thao
nnthao@fit.hcmus.edu.vn

Content outline

- Data objects and Attributes
- Basic statistical data descriptions
- Basic data visualization
- Data proximity measures
- Data correlation analysis

Exploratory data analysis (EDA)

- EDA analyzes data to summarize their main characteristics, using statistical graphics and data visualization methods.
- Data scientists use EDA to determine the best way to handle data sources to get the answers they need.



Initial Data Analysis (IDA)

- IDA focuses on examining the structure, quality, and basic characteristics of a dataset.
- Purposes:
 - Understand the basic structure and quality of the dataset.
 - Detect and fix problems (e.g., missing values, incorrect types, inconsistencies).
 - Ensure the data is clean, valid, and ready for further analysis.
- IDA makes sure that the data is in good shape, while EDA explores well-prepared data to uncover insights.

Common tasks in EDA (and IDA)

```
reason: integer (nullable = true)
reasonDetail: struct (nullable = true)
|-- a: string (nullable = true)
|-- b: array (nullable = true)
|   |-- element: struct (containsNull = true)
|   |   |-- b1: string (nullable = true)
|   |   |-- b2: array (nullable = true)
|   |       |-- element: struct (containsNull = true)
|   |           |-- b3: string (nullable = true)
|   |           |-- b4: integer (nullable = true)
|   |           |-- b5: integer (nullable = true)
```

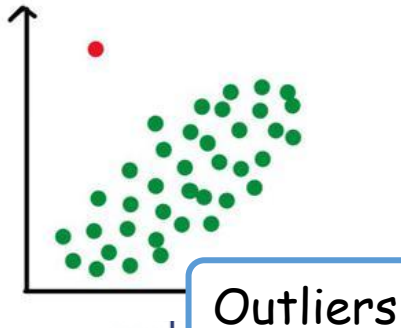
Data types

Flavor	Scoops sold	Contains chocolate?	Smooth or chunky?
Vanilla	300	No	Smooth
Chocolate	450	Yes	Smooth
Cookies & Cream	275	Yes	Chunky
Mint Chocolate Chip	315	Yes	Chunky

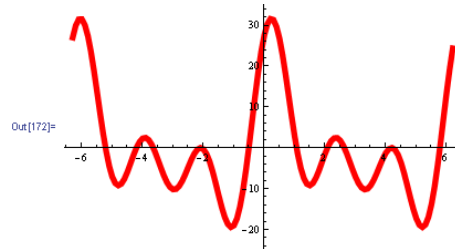
Data preprocessing



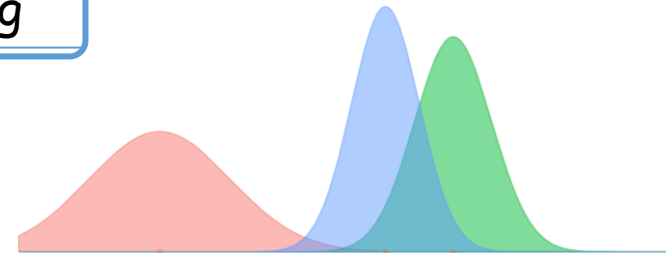
Data quality



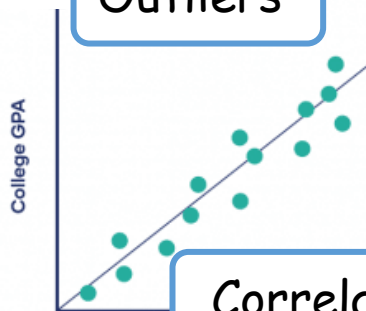
Outliers



Pattern discovery



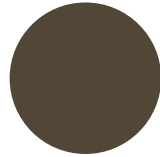
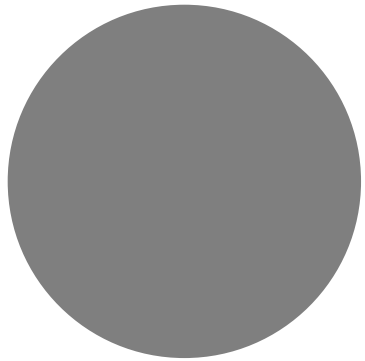
Data distribution



Correlation



Data visualization



Data objects and Attributes

Data collection: Record datasets

- Relational / transactional tuples
- Term-frequency vectors, numerical matrices, crosstabs

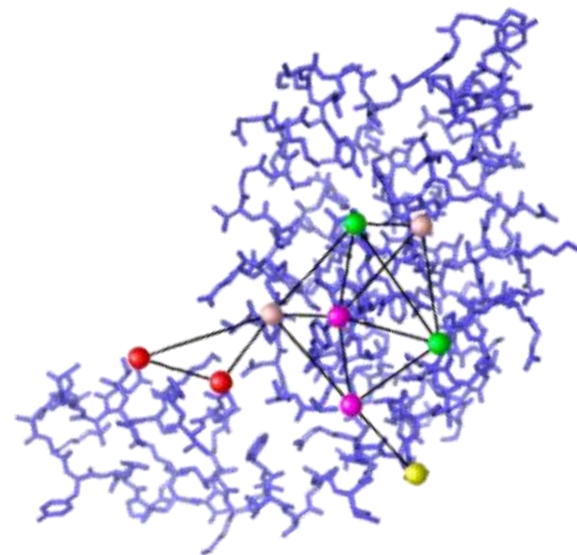
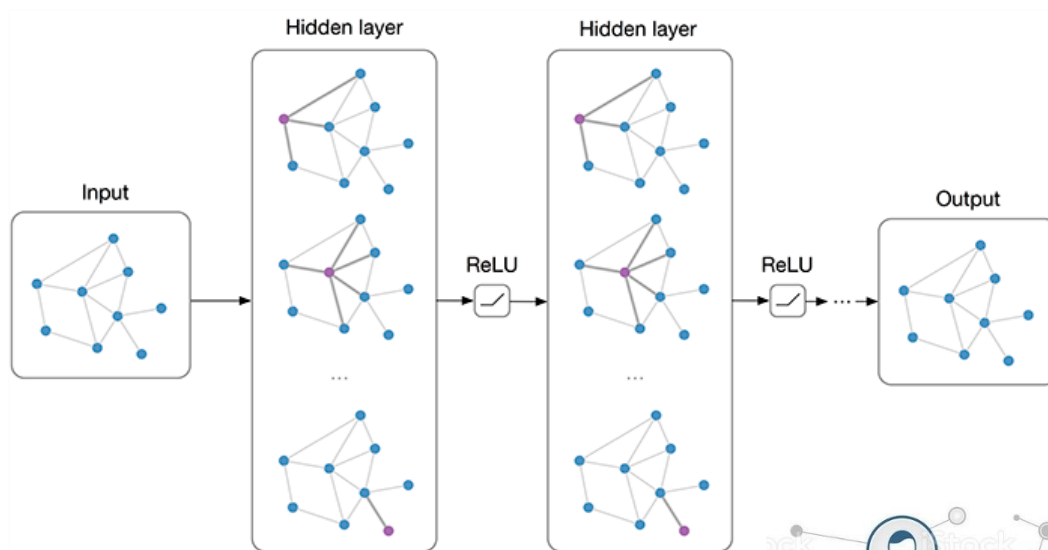
<i>TID</i>	Items
1	{Bread, Milk}
2	{Bread, Diapers, Beer, Eggs}
3	{Milk, Diapers, Beer, Cola}
4	{Bread, Milk, Diapers, Beer}
5	{Bread, Milk, Diapers, Cola}

	for	great	greatest	lasagna	life	love
sentence 1	0	0	1	0	1	1
sentence 2	0	2	0	0	0	1
sentence 3	0	0	1	0	0	1
sentence 4	1	0	0	1	0	1

		Task Performance		Total
		Fail	Success	
User Felt	Very bad	0	0	0
	Bad	2	1	3
	Neutral	1	4	5
	Good	0	15	15
	Very good	0	5	5
Total		3	25	28

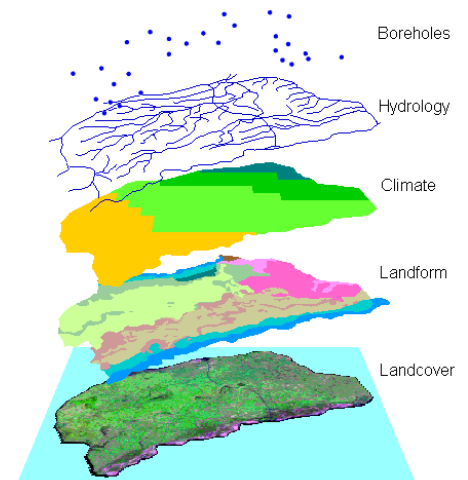
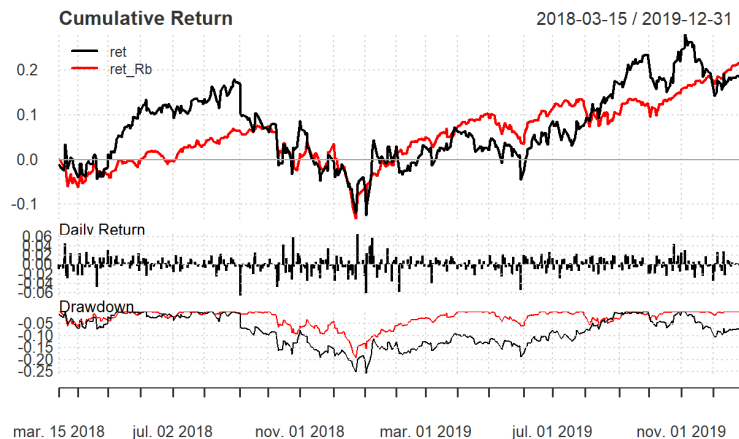
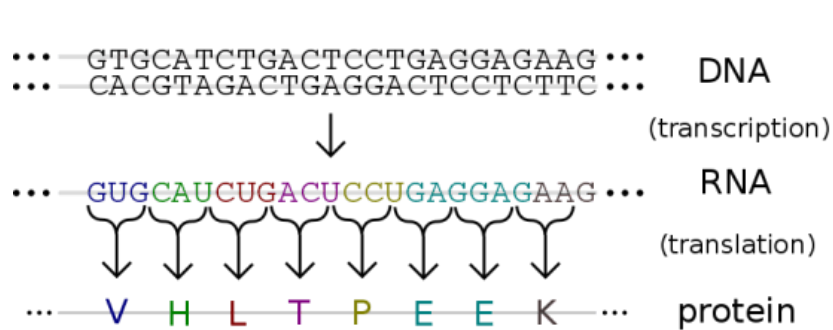
Data collection: Graph datasets

- The Internet, social networks, molecular structures



Data collection: Ordered datasets

- Sequential data: transaction sequences, genetic sequences
- Video data, temporal data, time-series data, etc.



Data objects

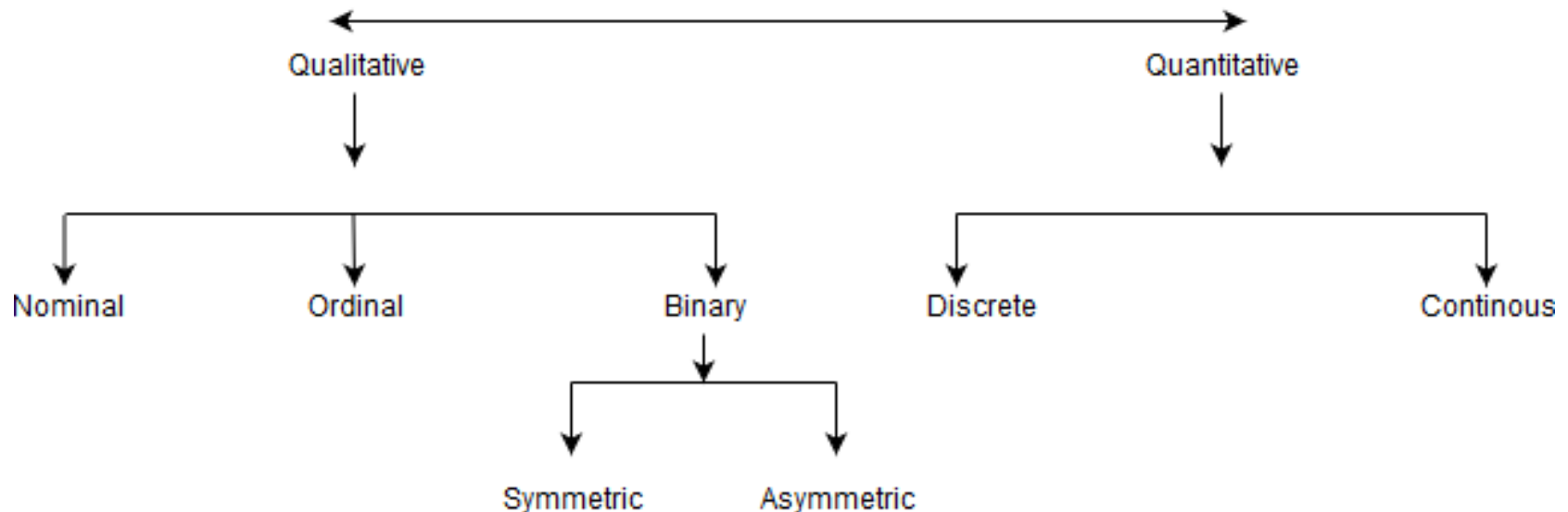
- A **data object** depicts an entity, serving as the **building block for a dataset**.
 - Similar terms: sample, example, instance, data point, and tuple



- Data objects are described by **attributes**.
 - In a database: rows → data objects, columns → attributes

Attributes

- An **attribute** shows some **characteristic of a data object**.
 - Similar terms: dimension, feature, and variable
 - E.g., a Customer object has 3 attributes {id, name, address}
- **Observation**: an observed value for a given attribute
- **Feature vector**: a set of attributes used to describe an object

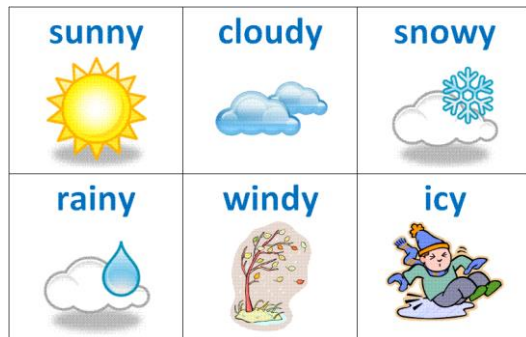


Attribute types: Nominal

- Qualitative, values do not have any meaningful order
- Enumerations: categories, states, or “names of things”



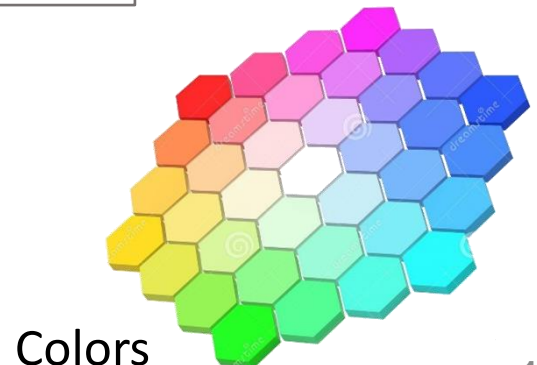
Day and Night



Weather



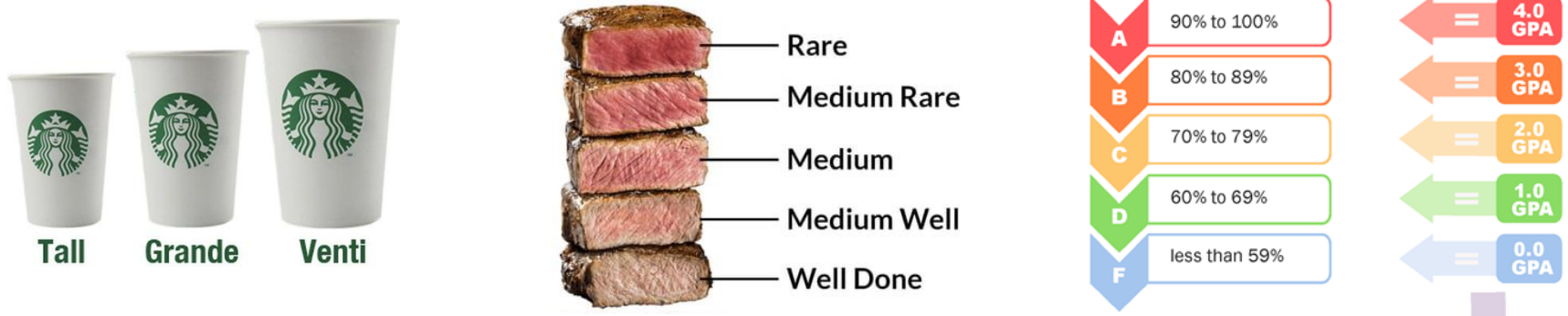
Occupation



Colors

Attribute types: Ordinal

- Qualitative, values have a meaningful order (ranking) but magnitude between successive values is not known



- Useful for subjective assessments of qualities that cannot be measured objectively
 - E.g., customer satisfaction

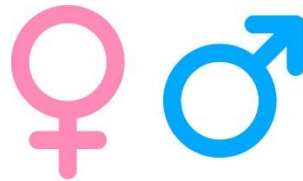


Attribute types: Binary

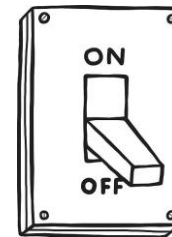
- Nominal attribute with only 2 states
- **Symmetric binary**: both outcomes **equally important**



Day and night

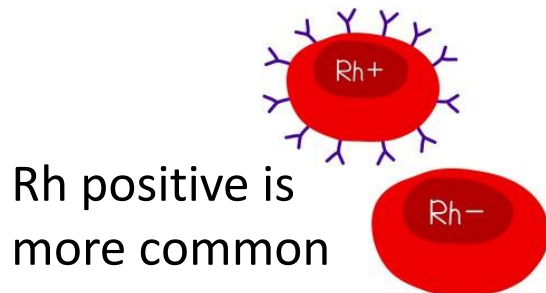


Male and Female



Switch light
On and Off

- **Asymmetric binary**: outcomes **not equally important**
 - Convention: assign 1 to the most important outcome (e.g., HIV test)



A positive result is more significant



Attribute types: Numeric

Interval numeric attribute

- Measured on a scale of **equal-sized units**
- Values have order (e.g., temperature in C° or F°, calendar dates)
- No true **zero-point**: able to compute the difference – not able to talk of one value as being a multiple of another
 - E.g., 20°C is five degrees higher than 15°C (right), 10°C is twice as warm as 5°C (wrong)

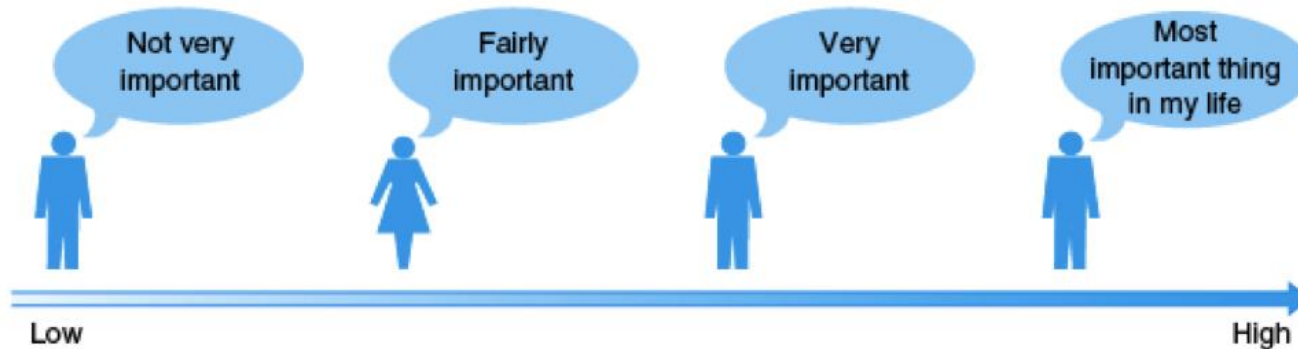
Ratio numeric attribute

- Inherent **zero-point**
- Values can be considered as being an order of magnitude larger than the unit of measurement
 - E.g., temperature (10°K is twice as high as 5°K), monetary (you are 100 times richer with \$100 than with \$1), measurements (height, weight)

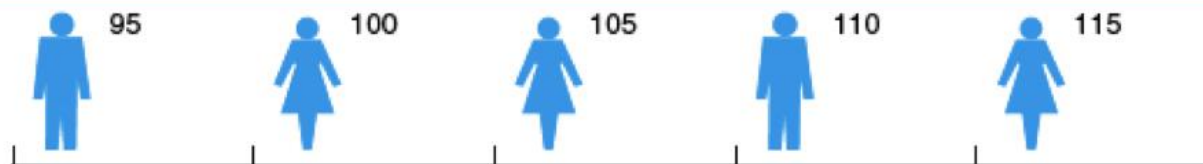
Nominal Measure Example: Gender



Ordinal Measure Example: Religiosity "How important is religion to you?"



Interval Measure Example: IQ



Ratio Measure Example: Income



Attributes: Discrete vs. Continuous

- There are many ways to organize attribute types, which are not mutually exclusive.
- Discrete attribute
 - Only a finite or countably infinite set of values
 - The values are sometimes represented as integers.
 - Binary attributes are a special case of discrete attributes.
- Continuous attribute
 - Real numbers of continuous domains
 - The values are usually represented using a finite number of digits
→ floating-point variables

Quiz 01: Data types

1. For each of the following pairs of data types, given an example to contrast the characteristics of data types.

- Nominal data vs. Ordinal data
- Symmetric binary data vs. Asymmetric binary data
- Interval numeric data vs. Ratio numeric data

2. How to check the data type of a variable in **Python**?

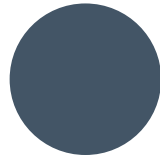
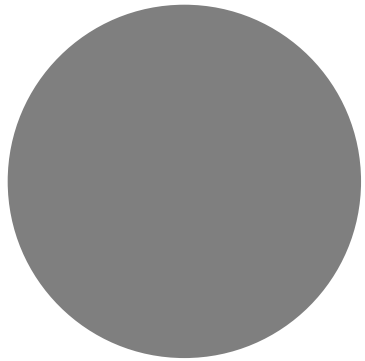
Show the data type of the three variables, x, y, z, shown aside.

```
x = 42  
y = "Hello"  
z = [1, 2, 3]
```

3. How to check the schema of a **pandas** Dataframe?

Check the scheme of the Dataframe show aside.

```
df = pd.DataFrame(  
    "name": ["Alice", "Bob"],  
    "age": [25, 30],  
    "salary": [50000.0, 60000.0]  
)
```



Basic statistical data descriptions

Central tendency: Arithmetic mean

- Let x_1, x_2, \dots, x_N be a set of N values or observations for some numeric attribute X .
- The **arithmetic mean** is defined as $\mu = \frac{1}{N} \sum_{i=1}^N x_i$
- The **weighted arithmetic mean** is written as $\mu^w = \frac{\sum_{i=1}^n w_i x_i}{\sum_{i=1}^n w_i}$
 - where w_i is the weight value that associates with x_i .
- It is the most common and effective numeric measure

Central tendency: Arithmetic mean

- Consider the score records of John and Kelly.
- The (non-weighted) mean scores are

$$\mu_{John} = 82.6, \quad \mu_{Kelly} = 84.6$$

John's record		Kelly's record	
Homework	92	Homework	100
Quiz	74	Quiz	82
Lab	83	Lab	95
Test	76	Test	70
Final exam	88	Final exam	76

Homework	15 %
Quiz	10 %
Lab	20 %
Test	25 %
Final exam	30 %

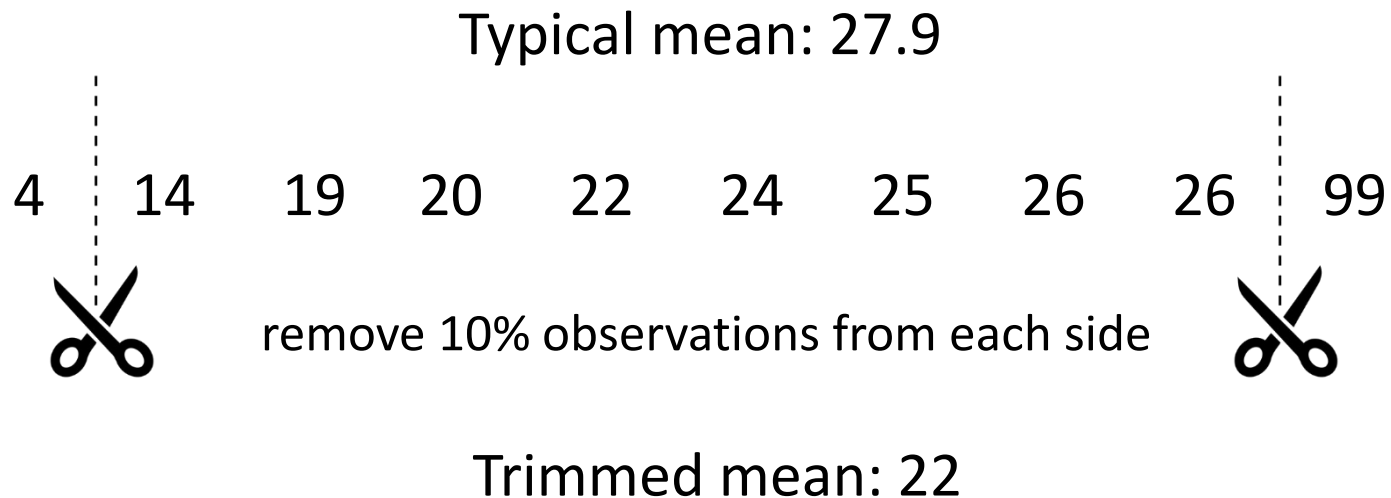
- We now have the course grade distribution.
- The weighted mean scores are

$$\mu_{John}^w = 83.2, \quad \mu_{Kelly}^w = 82.5$$

$$\mu_{John}^w = \frac{0.15 \times 92 + 0.1 \times 74 + 0.2 \times 83 + 0.25 \times 76 + 0.3 \times 88}{0.15 + 0.1 + 0.2 + 0.25 + 0.3} = 83.2$$

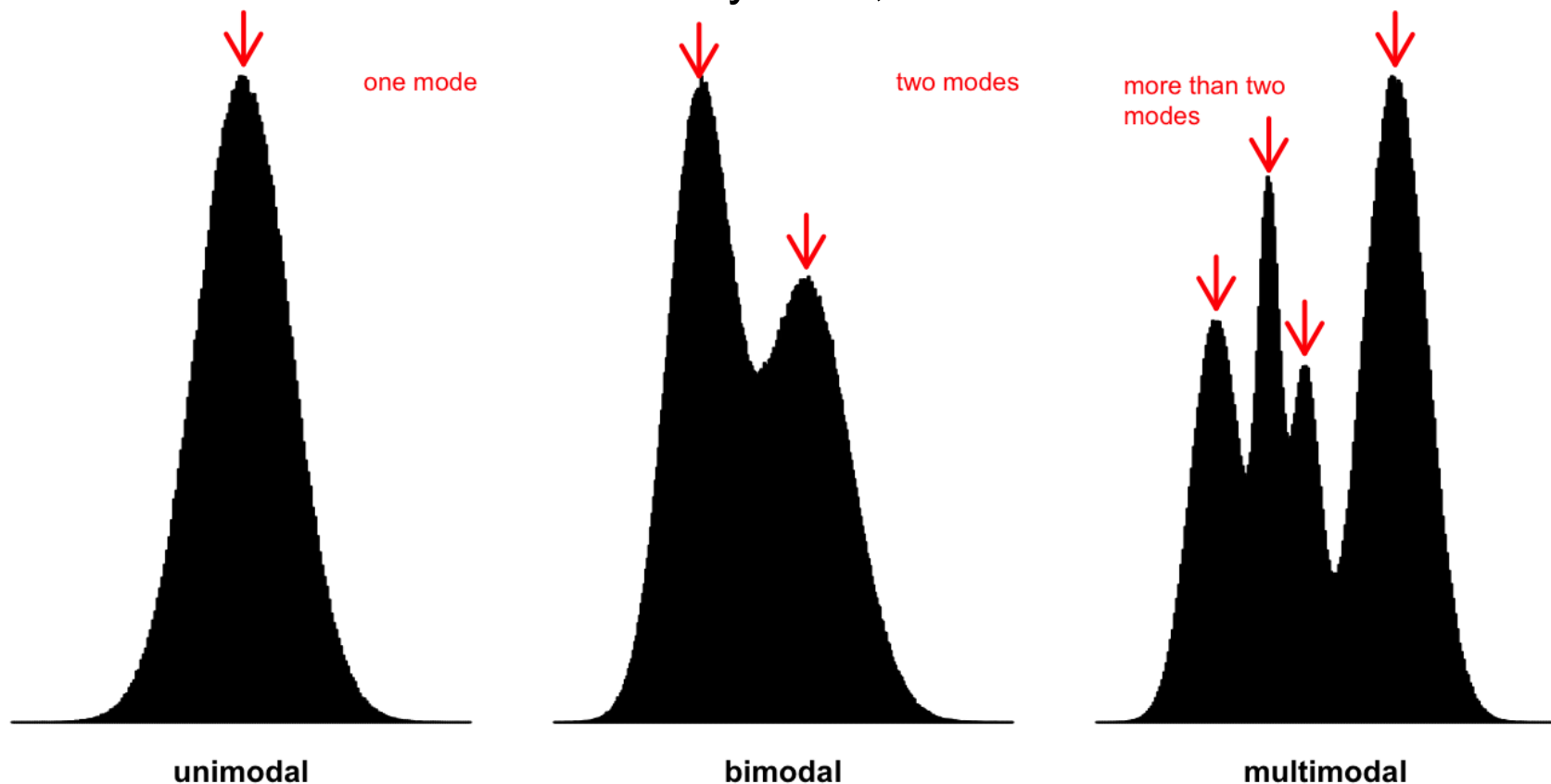
Central tendency: Arithmetic mean

- Means are highly sensitive to extreme values (e.g., outlier).
- Trimmed mean:** chop extreme values before calculating the regular mean



Central tendency: Mode

- **Mode** is the value that **occurs most frequently** in the data, defined for both qualitative and quantitative attributes.
 - If each data value occurs only once, then there is no mode



Central tendency: Median

- Suppose that the given set of N observations is sorted.
- **Median** is the **middle value** of the ordered set.
 - N is **odd**: pick the *exact middle value*; otherwise, take the *average of the two middlemost values*.
- **Midrange** is the **average of the largest and smallest values** in the set.

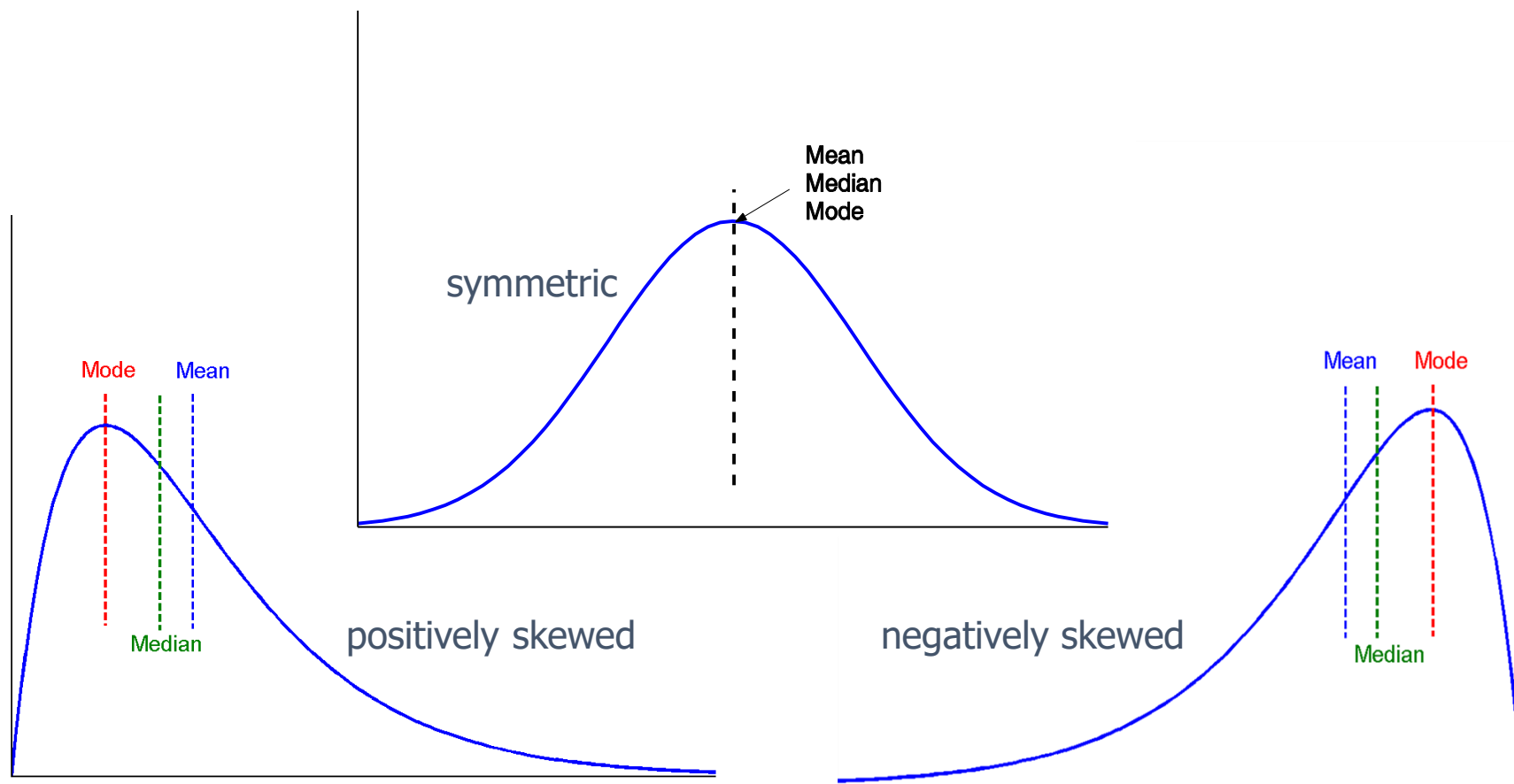
4 4 4 9 15 15 15 27 37 48

mean = 17.8 – mode: 4 and 15 – midrange = 26, median = $(15+15)/2 = 15$

3 3 6 9 15 15 15 27 27 37 48

mean = 18.636 – mode: 15 – midrange = 25.5, median = 15

Symmetric data vs. Skew data



- For moderately skewed unimodal numeric data, the empirical formula is
$$mean - mode \approx 3 \times (mean - median)$$

Quiz 02: Mean, mode, and Midrange

1. Consider the following 1D data series, which includes 13 data points.

31, 40, 19, 45, 5, 18, 30, 5, 33, 33, 25, 5, 20

Compute the following values: arithmetic mean, midrange, median, and mode.

2. For each of the following values, identify whether **pandas** provides a function to calculate the value.
 - Arithmetic mean
 - Median, mode
 - Weighted arithmetic mean
 - Midrange

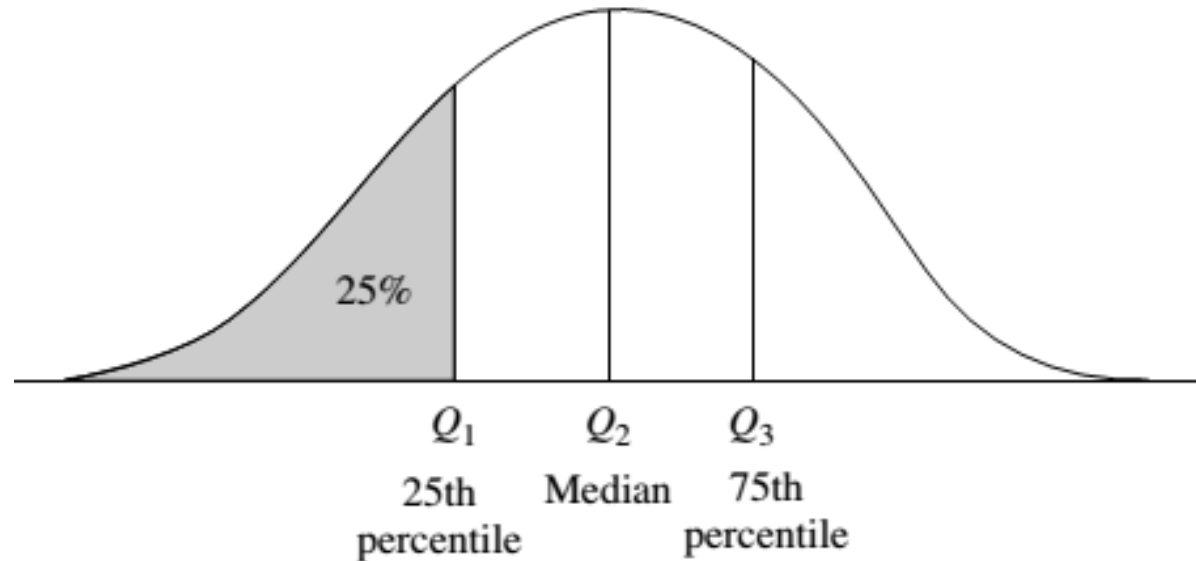
For each available function, show the result for the above data.

Data dispersion: Quantiles

- Let x_1, x_2, \dots, x_N be a set of N observations sorted in increasing order for a numeric attribute X .
- **Quantiles** are **points taken at regular intervals** of a data distribution, dividing it into equal-sized consecutive sets.
- **k^{th} q-quantile** ($0 < k < q, k \in \mathbb{N}^*$): a value x such that at most k/q data values $< x$ and at most $(q - k)/q$ of which $> x$.
 - There are $q - 1$ q-quantiles.

Data dispersion: Quantiles

- **Quartiles** (4-quantiles) split the data distribution into four equal parts.



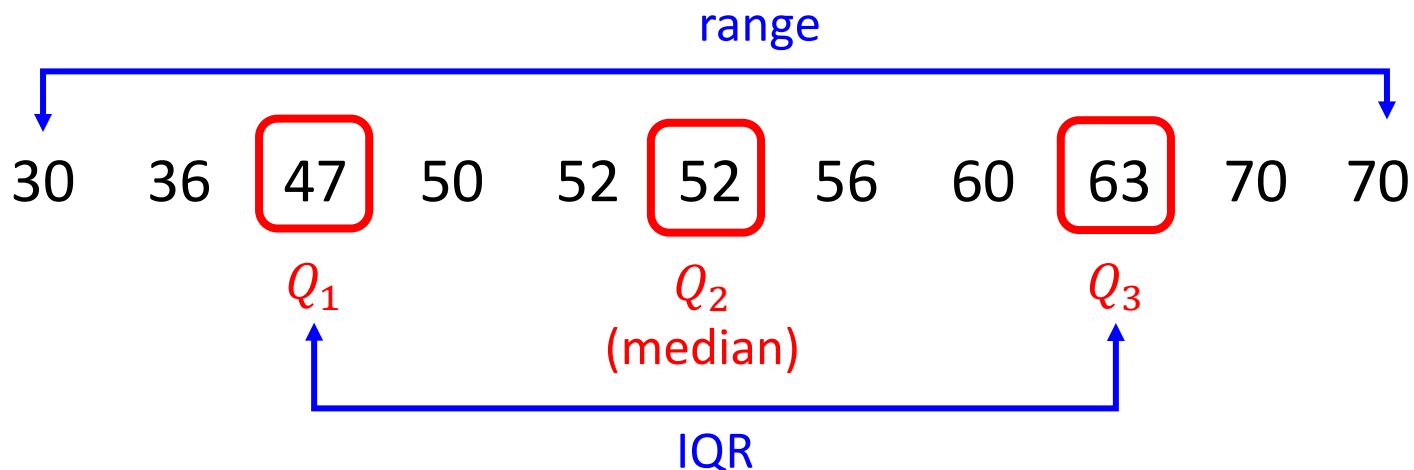
- **Percentiles** (100-quantiles): 100 equal-sized consecutive sets.
- **2-quantile** is the median that splits the distribution into halves.

Data dispersion: Interquartile range

- **Interquartile range** (IQR) is the distance between the first and third quartiles.

$$IQR = Q_3 - Q_1$$

- **Range** is the difference between the largest and smallest values in the set.



How to determine the quartile?

- Use the median to divide the ordered set into two halves.
 - If the original set has an even number of points, split it exactly in half
 - Otherwise, **do not include** the median in either half.
- Q_1 and Q_3 are the medians of the lower and upper halves, respectively.

6 7 15 36 39 40 41 42 43 47 49

Q_1 Q_2 Q_3

7 15 36 39 40 41

Q_1 $Q_2 = 37.5$ Q_3

Quiz 03: Quantiles

1. You are given the following dataset representing the scores of 15 students in a math exam, already sorted.

45, 48, 52, 55, 62, 67, 70, 72, 75, 77, 80, 85, 87, 90, 95

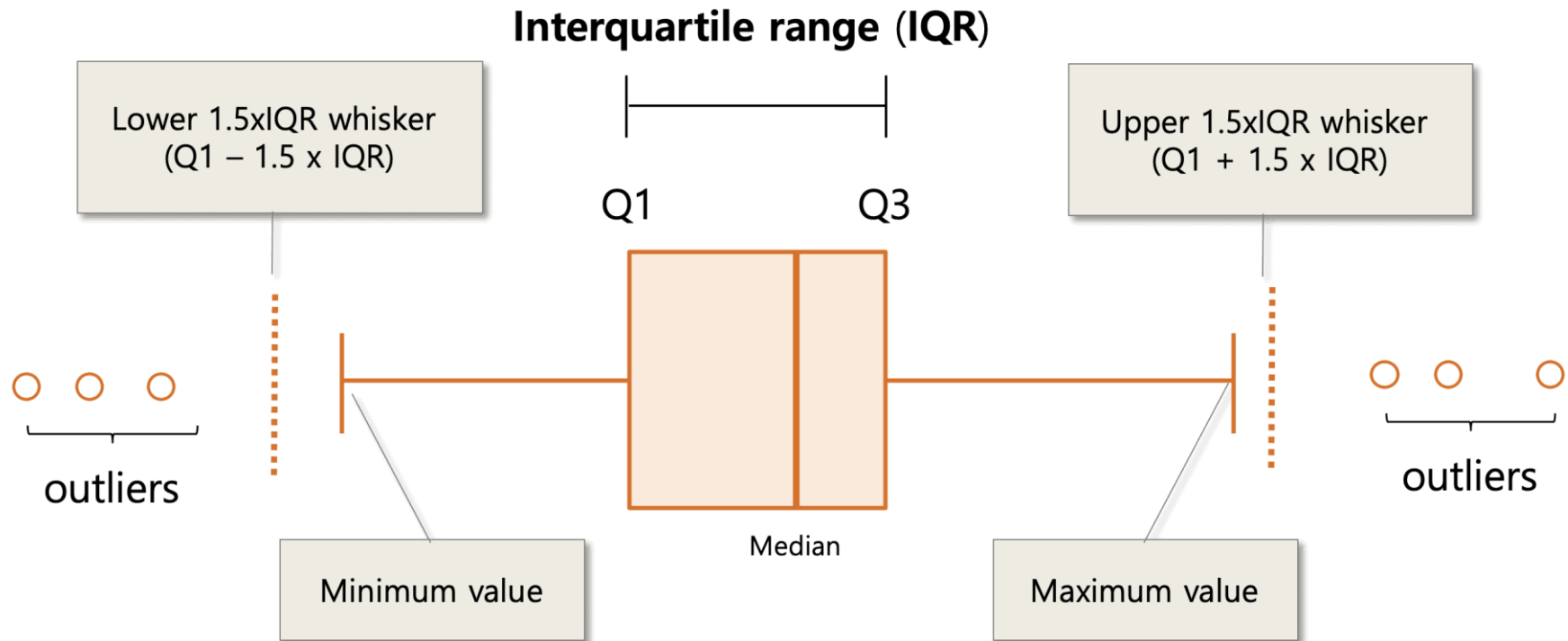
Compute the first quartile (Q1), second quartile (Q2), and third quartile (Q3), and IQR.

2. Identify whether **pandas** supports the calculation of quartiles and IQR. For each available function, show the result for the above data.

Data dispersion: Boxplot

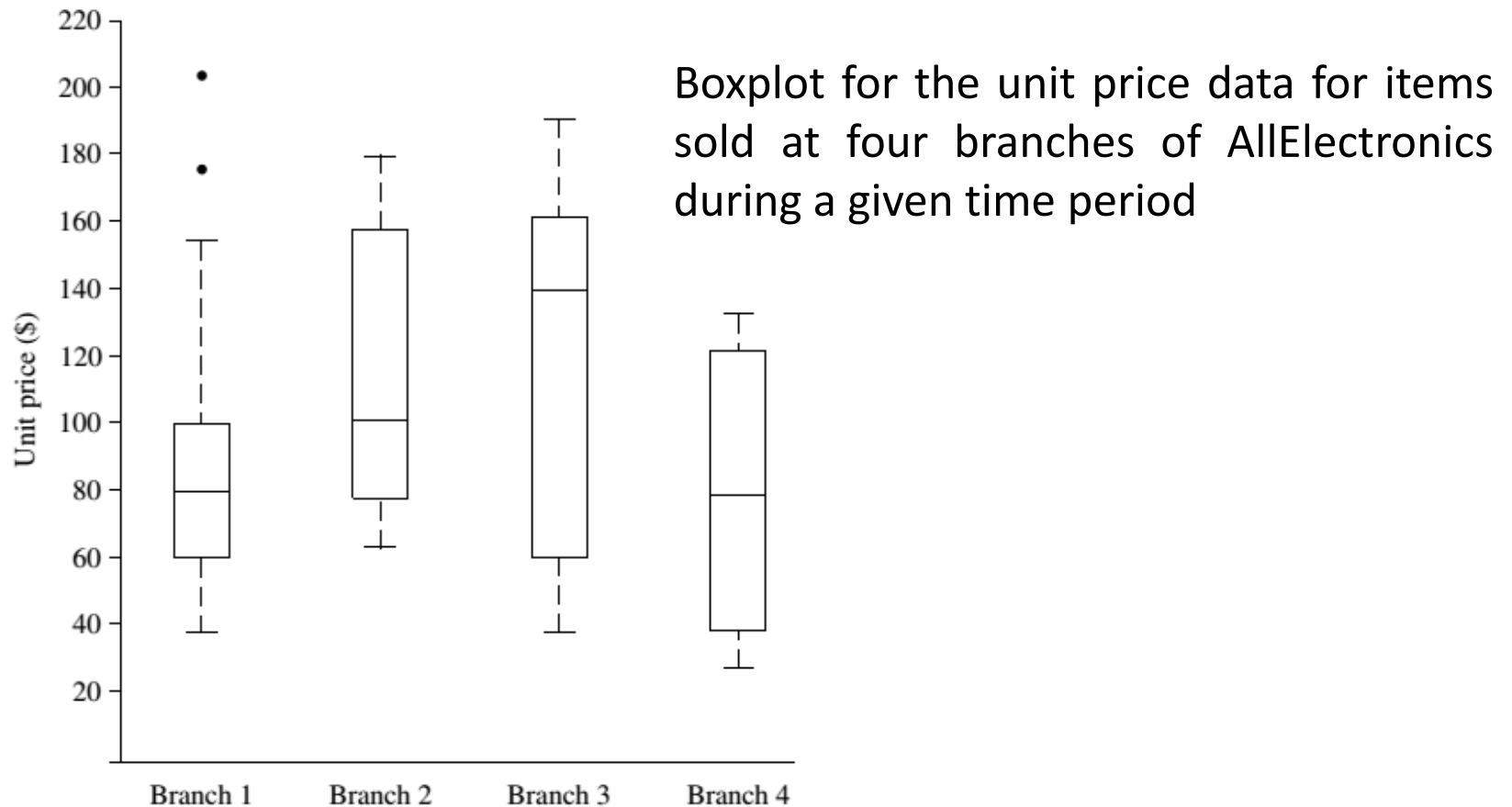
- A **five-number summary** of a distribution includes
 - The median (Q_2), the quartiles Q_1 and Q_3 ,
 - The smallest (*Min*) and largest (*Max*) individual values.
- This summary is presented by a **boxplot**.

Data dispersion: Boxplot



- The two whiskers refers to the smallest and largest values within $[Q_1 - 1.5 \times IQR, Q_3 + 1.5 \times IQR]$.
- Outliers: points that are out the above range, plotted individually

Data dispersion: Boxplot



- For Branch 1, the median price of items sold is \$80, Q_1 is \$60, and Q_3 is \$100. Notice that two outlying observations, 175 and 202, were plotted individually as they are more than $1.5 \times \text{IQR}$.

Quiz 04: Draw a box plot

1. Consider the following 1D data series, which includes 15 data points sorted in ascending order.

21, 25, 27, 29, 32, 36, 36, 48, 67, 70, 74, 75,
79, 150, 197

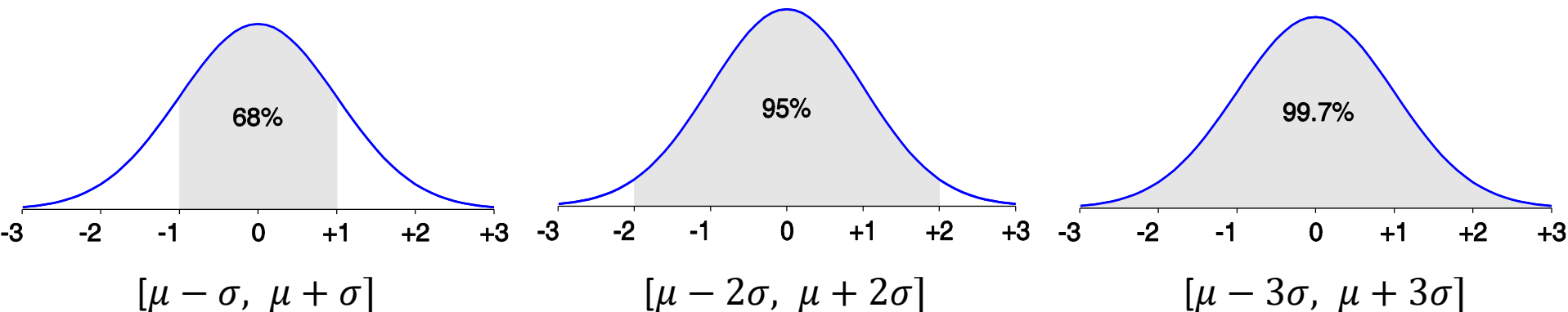
- Define the five-number summary for the above data.
 - Draw the boxplot representing the above five-number summary. Note the vertical axis and all the values.
2. How to draw a boxplot in **Python**?
 3. Can **scikit-learn** be used to draw a **boxplot**? If yes, how?

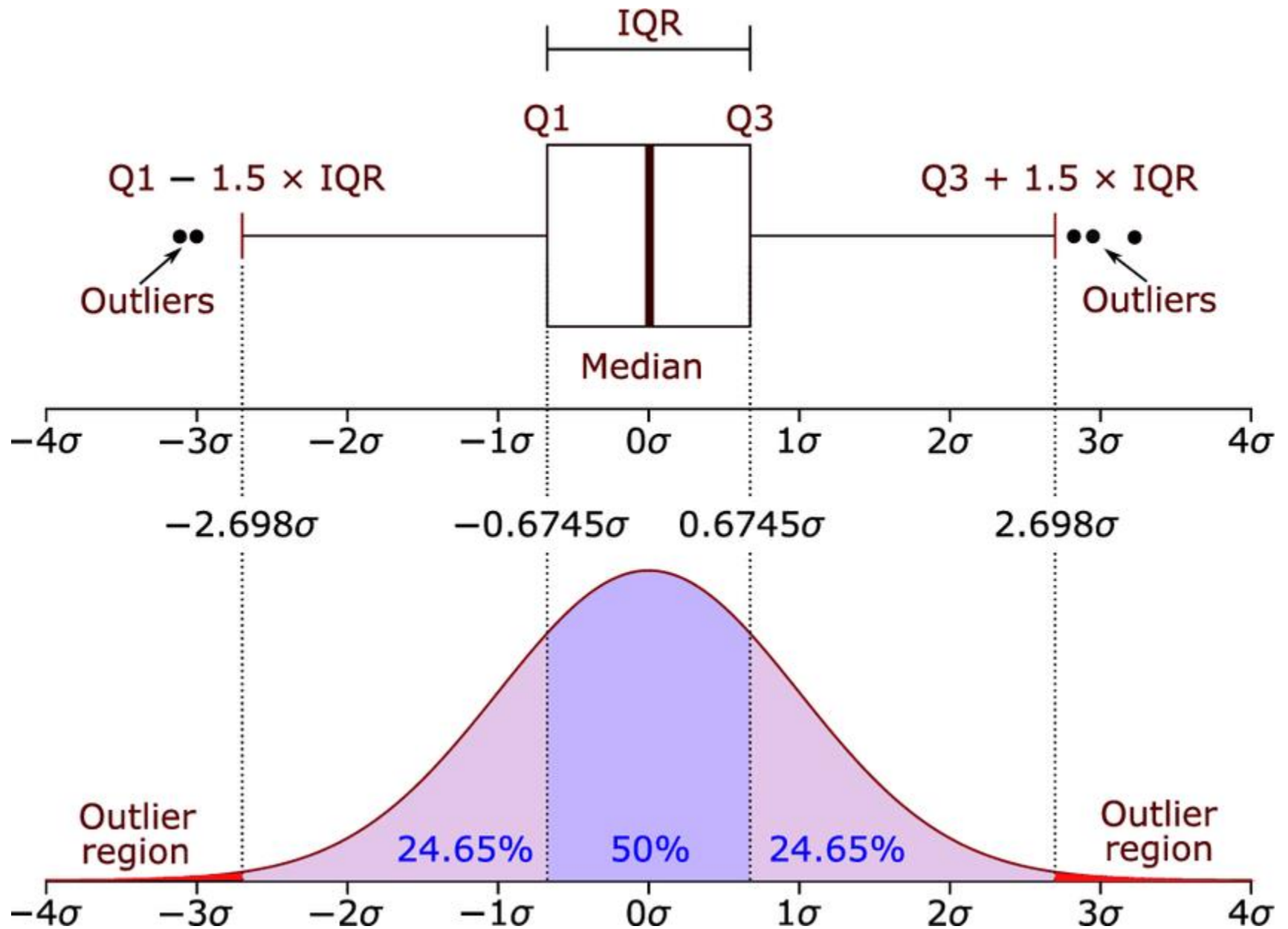
Data dispersion: Variance

- The (population) **variance** is defined as

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2 = \left(\frac{1}{N} \sum_{i=1}^N x_i^2 \right) - \bar{x}^2$$

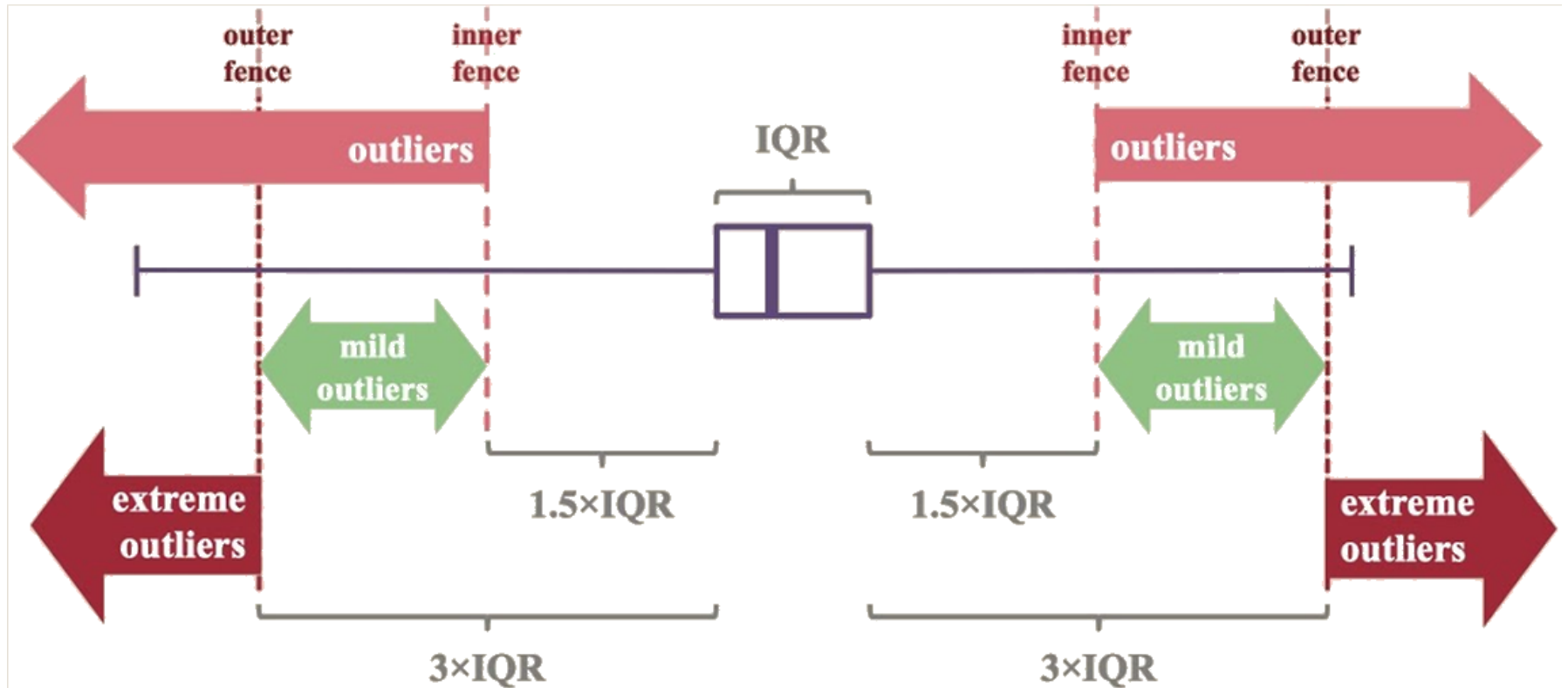
- The **standard deviation** is the square root of the variance.
 - Low $\sigma \rightarrow$ the data tends to be very close to the mean.*
 - High $\sigma \rightarrow$ the data spreads out over a large range of values.*





Box plot and probability density function of a normal distribution.

Types of outliers



Quiz 05: Variance and standard deviation

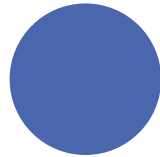
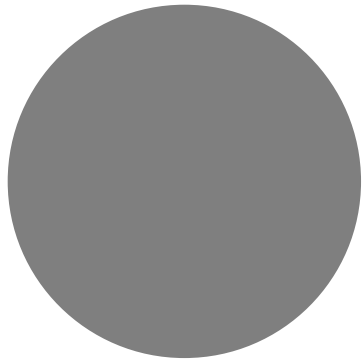
1. Consider the following 1D data series, which includes 15 data points sorted in ascending order.

21, 25, 27, 29, 32, 36, 36, 48, 67, 80, 84,
85, 89, 92, 97

Compute the variance and standard deviation.

2. Does **pandas** provide the variance and standard deviation for a dataframe? If yes, how?

For each available function, show the result for the above data.



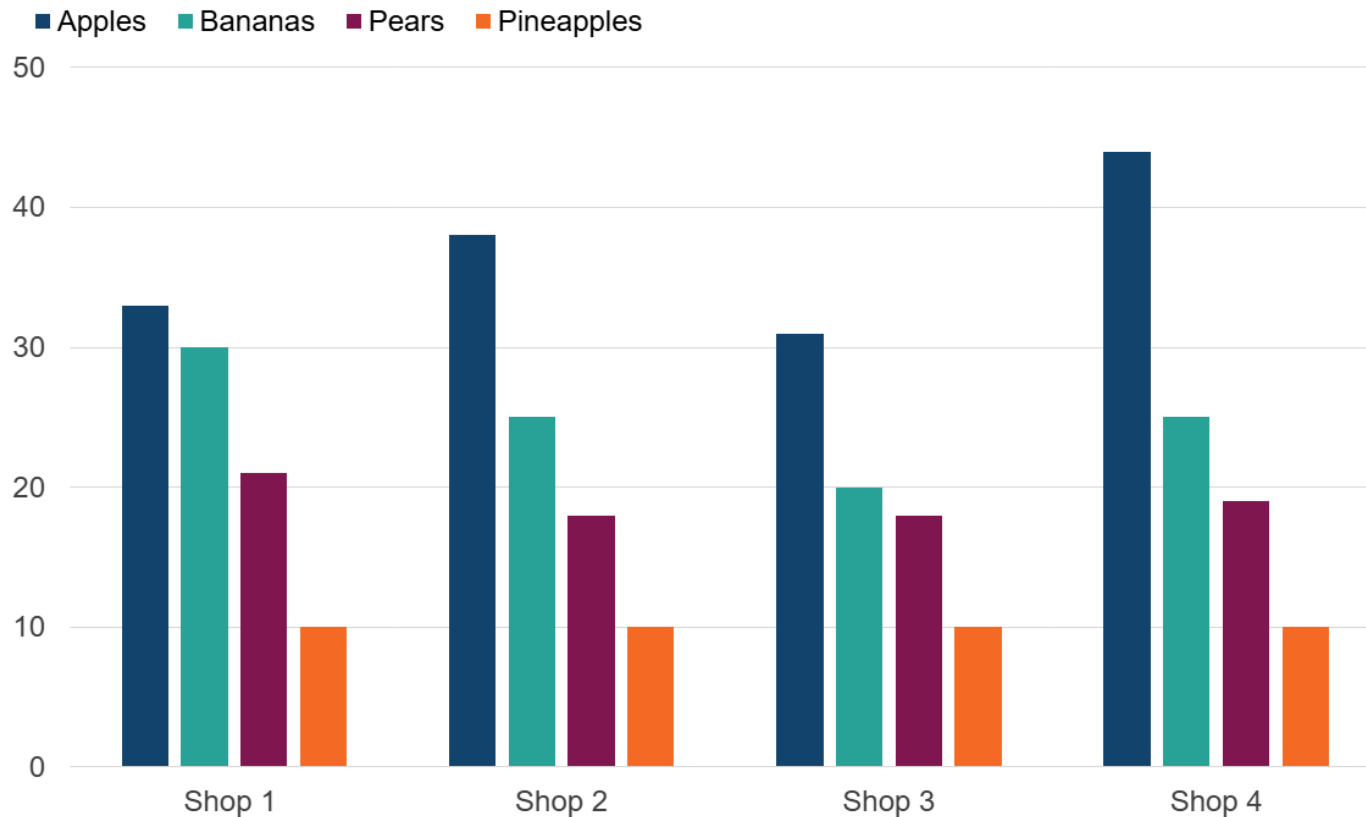
Basic data visualization

Why data visualization?

- Gain insight into an information space by mapping data onto graphical primitives
- Provide qualitative overview of large datasets
- Search for patterns, trends, irregularities, relationships
- Help find interesting regions and suitable parameters for further quantitative analysis
- Provide a visual proof of computer representations derived

Bar chart

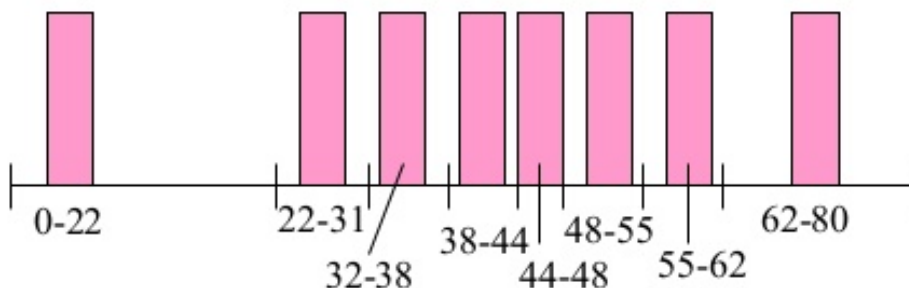
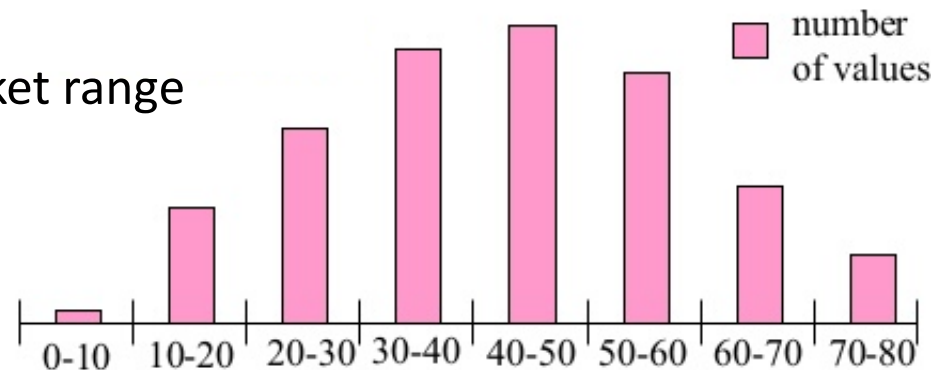
- A **bar chart** presents **nominal data** by using rectangular bars with heights proportional to the values represented.



Histogram

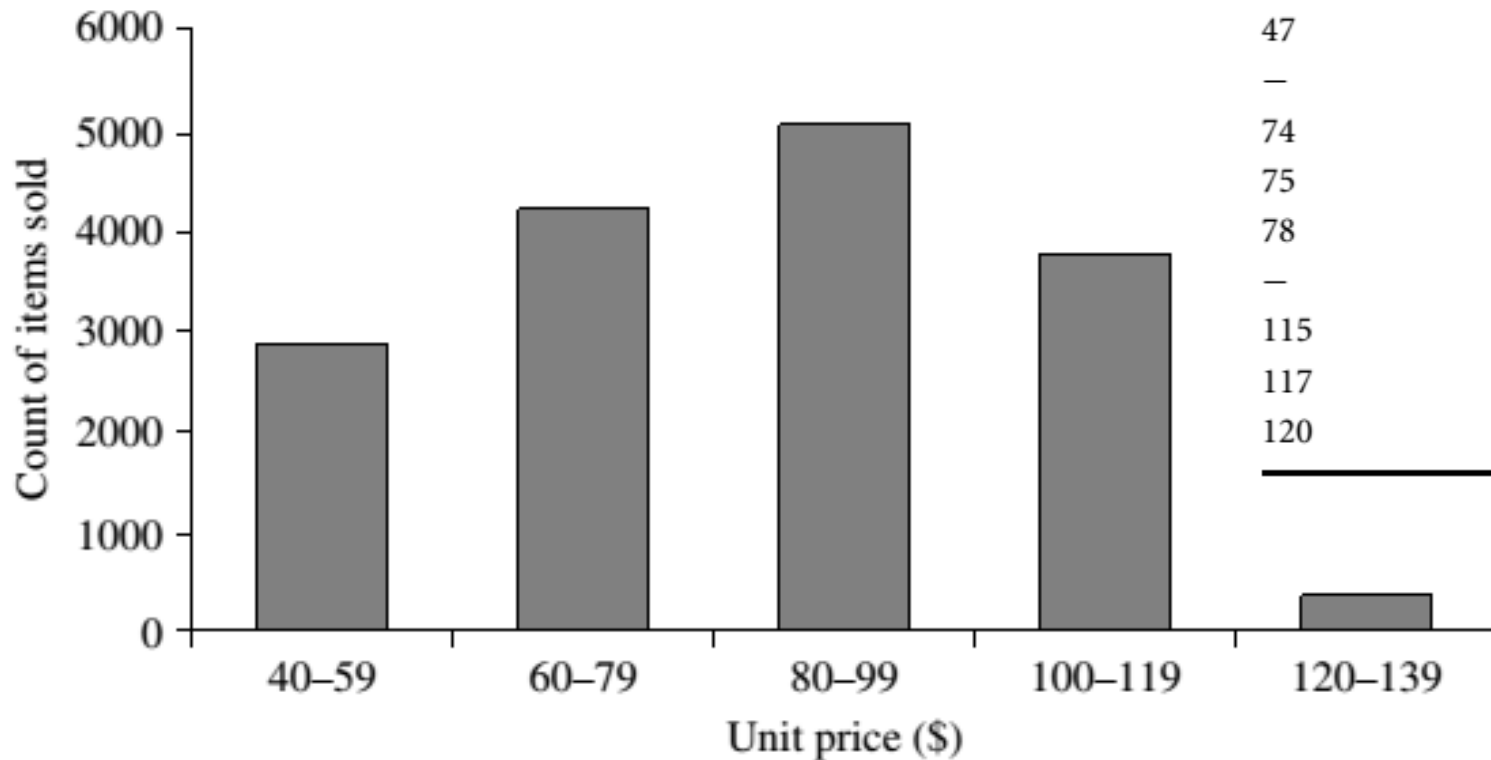
- The **range of values** for a **numeric attribute** X is **partitioned into disjoint consecutive subranges**, called **buckets** or **bins**.
- Each bar is for a subrange such that **its height represents the total items within the subrange**.

Equal-width: equal bucket range



Equal-frequency: equal bucket depth

Histogram: An example

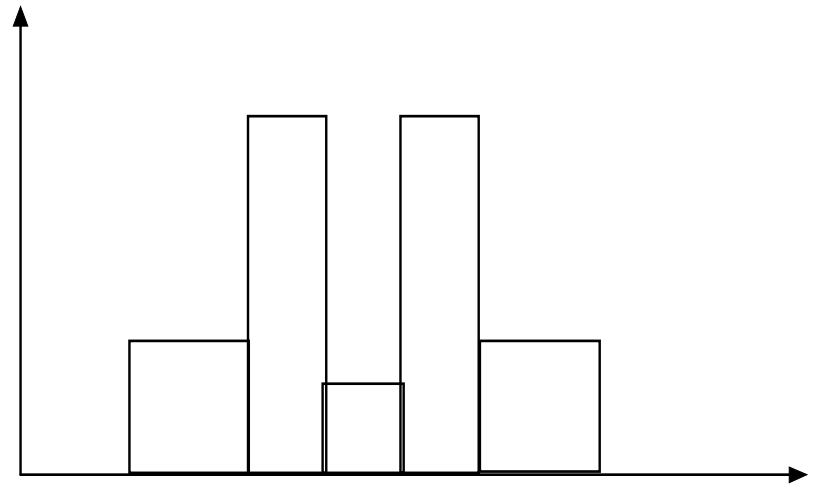
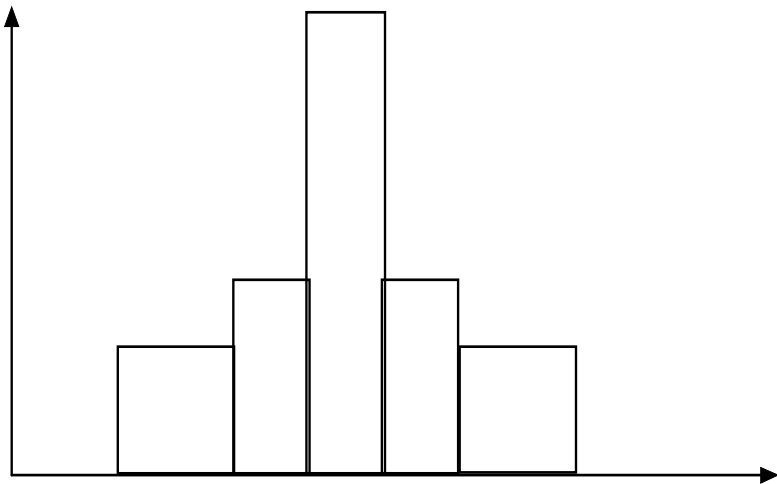


A Set of Unit Price Data for Items Sold at a Branch of *AllElectronics*

Unit price (\$)	Count of items sold
40	275
43	300
47	250
—	—
74	360
75	515
78	540
—	—
115	320
117	270
120	350

Histogram over boxplot

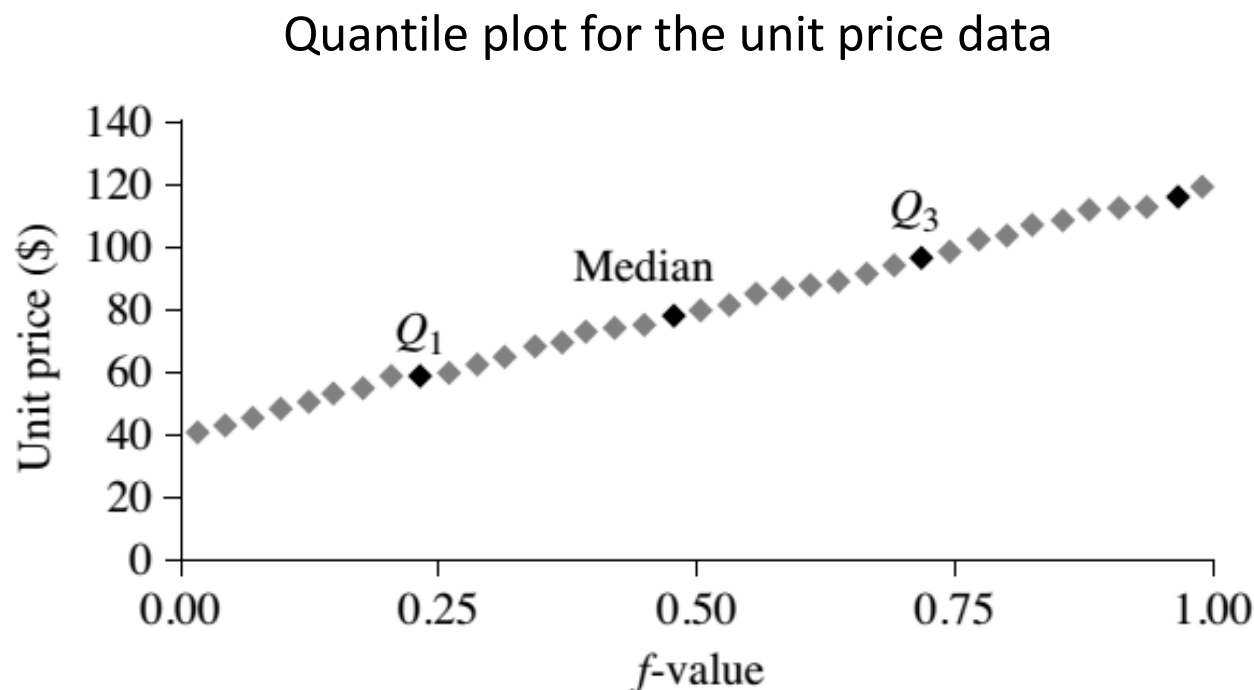
- The two following histograms may have the same boxplot.
- However, they represent rather different data distributions.



Quantile plot

- A **quantile plot** presents the plot quantile information for a **univariate data distribution**.
 - It allows access to both overall behavior and unusual occurrences.
- Let x_1, x_2, \dots, x_N be the data observations sorted in increasing order for some ordinal or numeric attribute X .
- Each value x_i is paired with $f_i = \frac{i-0.5}{N}$, indicating that approximately $f_i \times 100\%$ of data are $\leq x_i$.

Quantile plot: An example



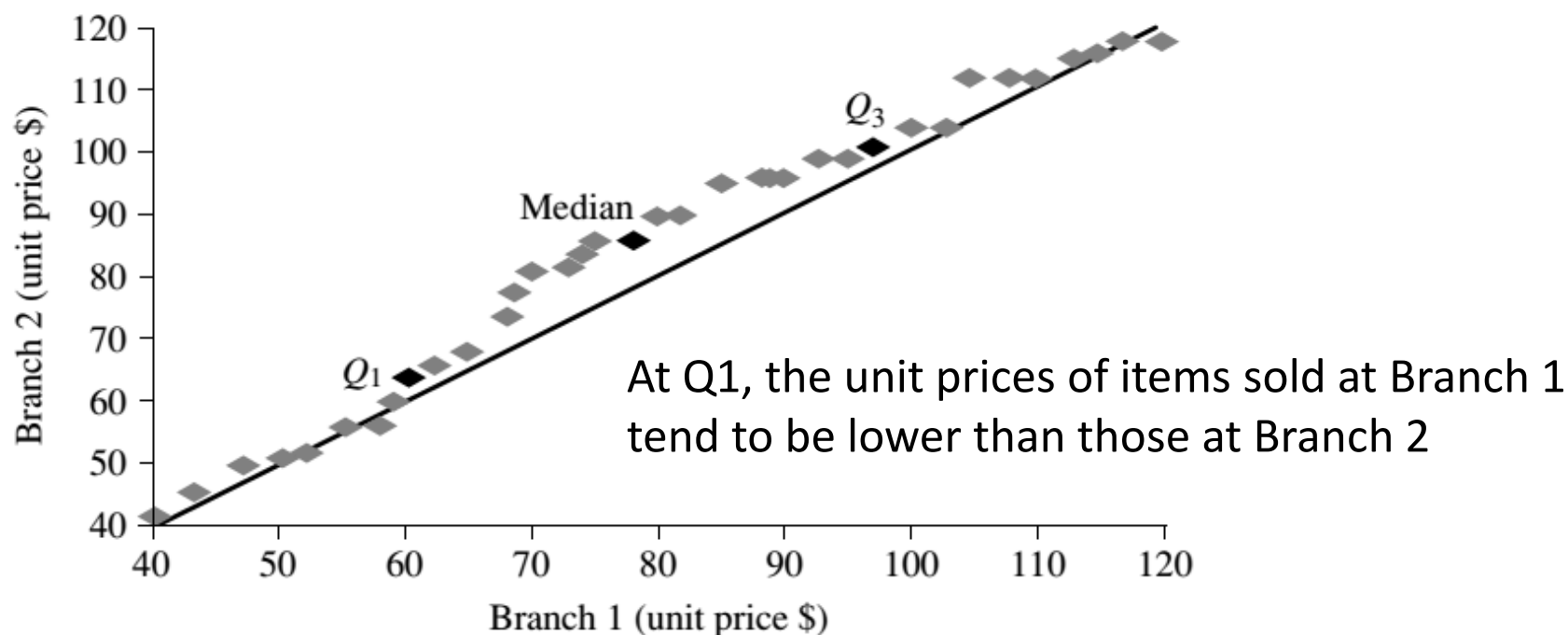
A Set of Unit Price Data for Items Sold at a Branch of *AllElectronics*

Unit price (\$)	Count of items sold
40	275
43	300
47	250
—	—
74	360
75	515
78	540
—	—
115	320
117	270
120	350

Quantile-Quantile plot

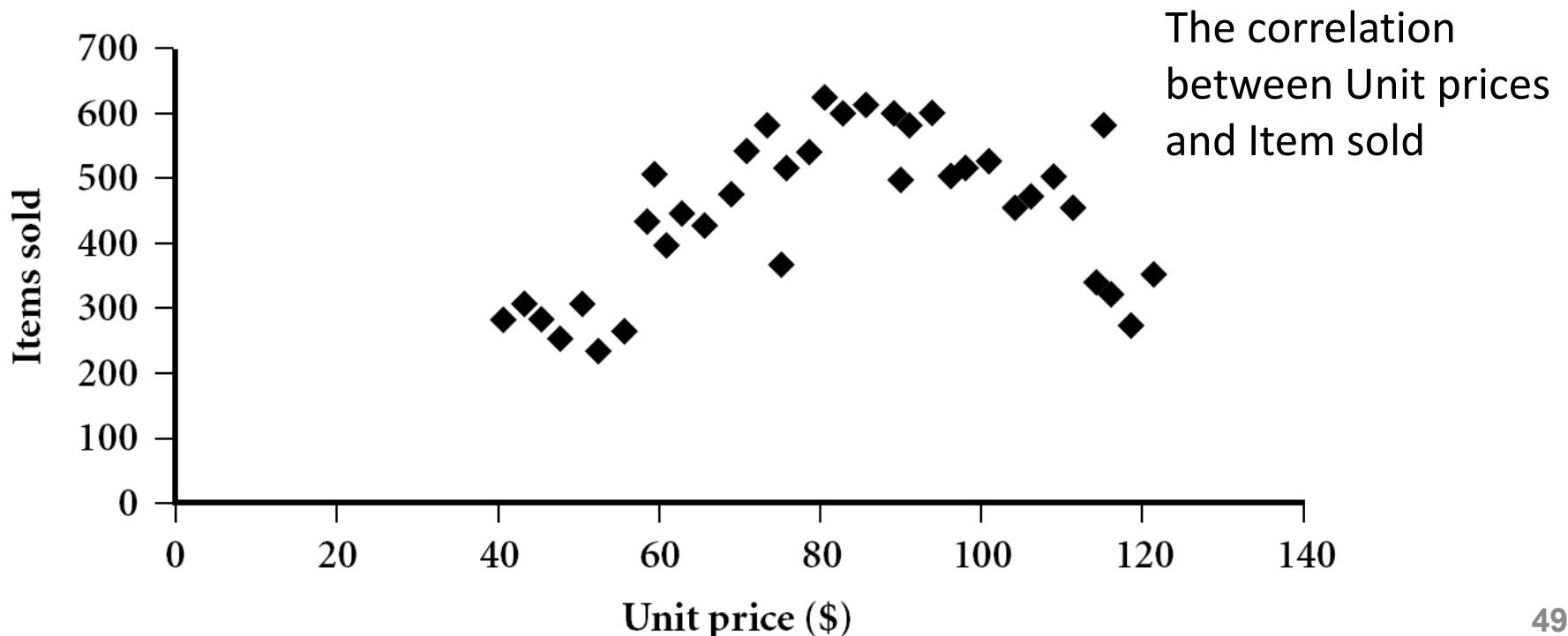
- A **quantile-quantile plot** draws the **quantiles of one univariate distribution** against the corresponding **quantiles of another**.

Is there a shift in going from one distribution to another?

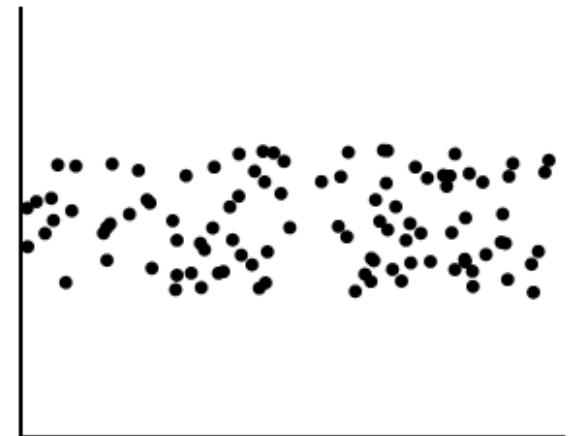
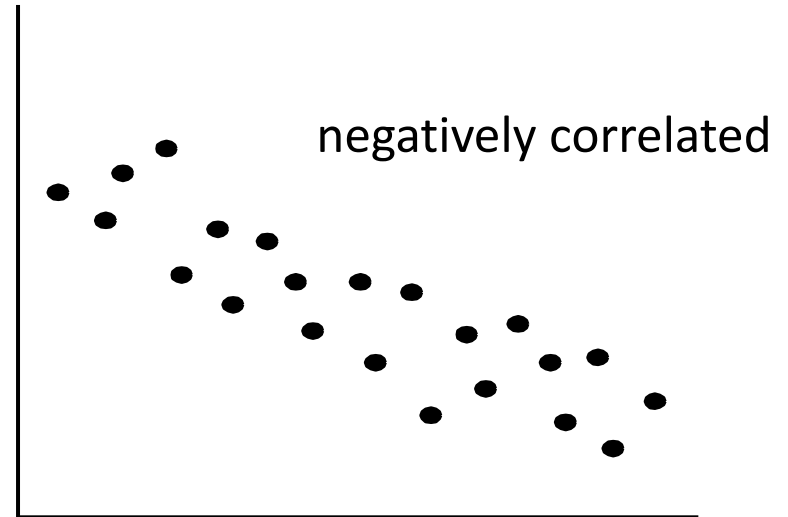
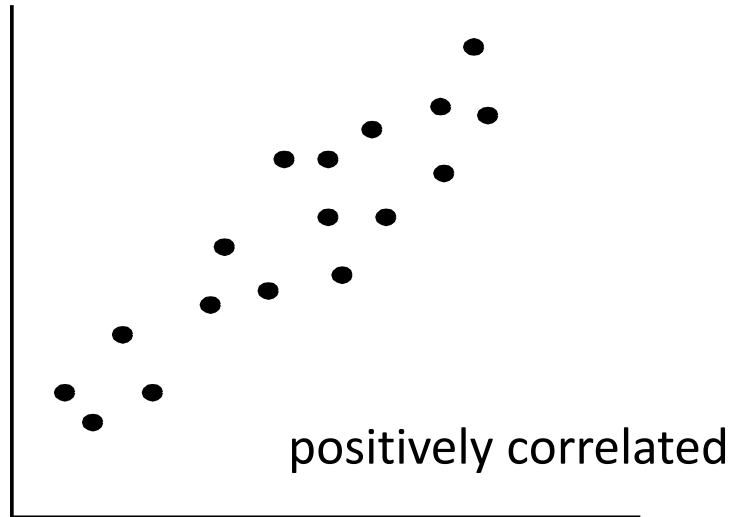


Scatter plot

- A **scatter plot** looks at the **bivariate data** to see clusters of points or outliers
 - Each pair of values is treated as a pair of coordinates and plotted as points in the plane.



Scatter plot: Data correlation



uncorrelated data

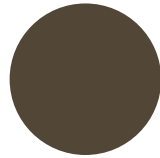
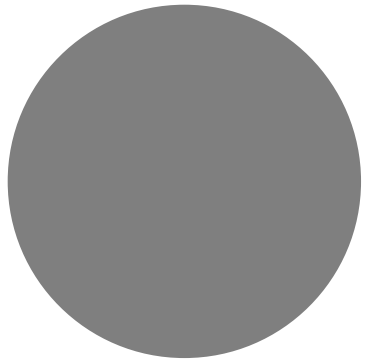
Quiz 06: Scatter plot

1. Consider the following data table, in which there are five tuples of two attributes, A and B.

No.	Attributes	
	A	B
1	19	16
2	25	10
3	13	26
4	12	29
5	16	20

Draw the scatter plot, whose the horizontal axis denotes attribute A, and the vertical axis represents attribute B.

2. How to draw a scatter plot using some **Python library**? Draw the scatter plot for the above data.



Data proximity measures

Similarity and Dissimilarity

Similarity

- A numerical measure of **how alike** two data objects, i and j , are
- Values often falls in the range $[0,1]$: 0 – unlike \rightarrow 1 – identical

Dissimilarity (distance)

- A numerical measure of **how different** two data objects are
- It works in an opposite direction to some similarity measure
- The lower bound is often 0, while the upper limit varies

Proximity

- This refers to either similarity or dissimilarity

Feature matrix vs. Dissimilarity matrix

- Feature matrices are essential to most machine learning task.

Feature matrix

$$\begin{bmatrix} x_{11} & \dots & x_{1f} & \dots & x_{1p} \\ \dots & \dots & \dots & \dots & \dots \\ x_{i1} & \dots & x_{if} & \dots & x_{ip} \\ \dots & \dots & \dots & \dots & \dots \\ x_{n1} & \dots & x_{nf} & \dots & x_{np} \end{bmatrix}$$

- n data points with p dimensions
- Object-by-attribute structure

Dissimilarity matrix

$$\begin{bmatrix} 0 & & & & \\ d(2,1) & 0 & & & \\ d(3,1) & d(3,2) & 0 & & \\ \vdots & \vdots & \vdots & & \\ d(n,1) & d(n,2) & \dots & \dots & 0 \end{bmatrix}$$

- A collection of distances for all pairs of n objects
- Object-by-object structure

- Many nearest-neighbor algorithms use dissimilarity matrices.

Measures for nominal attributes

- Let the number of states of a nominal attribute be M
- Method 1:** Simple matching $d(i, j) = \frac{p-m}{p}$
 - m : the number of attributes for which i and j are in the same state,
 - p : the total number of attributes describing the objects

Object Identifier	test-1 (nominal)	test-2 (ordinal)	test-3 (numeric)
1	code A	excellent	45
2	code B	fair	22
3	code C	good	64
4	code A	excellent	28

$$\begin{bmatrix} 0 & & & \\ 1 & 0 & & \\ 1 & 1 & 0 & \\ 0 & 1 & 1 & 0 \end{bmatrix}$$

$$\begin{bmatrix} d(2, 1) & & & \\ d(3, 1) & d(3, 2) & & 0 \\ d(4, 1) & d(4, 2) & d(4, 3) & 0 \end{bmatrix}$$

- Method 2:** Create a binary attribute for each of the M states
- Measures of similarity $sim(i, j) = 1 - d(i, j) = \frac{m}{p}$

Measures for binary attributes

- Contingency table

		Object j		
		1	0	sum
Object i	1	q	r	$q + r$
	0	s	t	$s + t$
	sum	$q + s$	$r + t$	p

- Symmetric binary variable

$$d(i, j) = \frac{r + s}{q + r + s + t}$$

- Asymmetric binary variable

$$d(i, j) = \frac{r + s}{q + r + s}$$

- Jaccard coefficient: $\text{sim}(i, j) = 1 - d(i, j) = \frac{q}{q + r + s}$

Measures for binary attributes

<i>name</i>	<i>gender</i>	<i>fever</i>	<i>cough</i>	<i>test-1</i>	<i>test-2</i>	<i>test-3</i>	<i>test-4</i>
Jack	M	Y	N	P	N	N	N
Jim	M	Y	Y	N	N	N	N
Mary	F	Y	N	P	N	P	N
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮

- Gender is symmetric binary, the remaining attributes are asymmetric
- Let the values Y and P be 1 and the value N be 0.
- Suppose that the distance between objects (patients) is computed based only on the asymmetric attributes

$$d(\text{Jack}, \text{Jim}) = \frac{1+1}{1+1+1} = 0.67, \quad d(\text{Jack}, \text{Mary}) = \frac{0+1}{2+0+1} = 0.33$$

$$d(\text{Jim}, \text{Mary}) = \frac{1+2}{1+1+2} = 0.75$$

Measures for numeric attributes

- Consider two data points of p -dimensional

$$i = (x_{i1}, x_{i2}, \dots, x_{ip}) \text{ and } j = (x_{j1}, x_{j2}, \dots, x_{jp})$$

- Minkowski distance (L_h norm)

$$d(i, j) = \sqrt[h]{|x_{i1} - x_{j1}|^h + |x_{i2} - x_{j2}|^h + \dots + |x_{ip} - x_{jp}|^h}$$

- where h is the order

Measures for numeric attributes

- $h = 1$: Manhattan (city block, L_1 norm) distance

$$d(i, j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \cdots + |x_{ip} - x_{jp}|$$

- $h = 2$: Euclidean (L_2 norm) distance

$$d(i, j) = \sqrt{|x_{i1} - x_{j1}|^2 + |x_{i2} - x_{j2}|^2 + \cdots + |x_{ip} - x_{jp}|^2}$$

- $h \rightarrow \infty$: “supremum” (L_{max} / L_∞ norm, Chebyshev) distance

$$d(i, j) = \lim_{h \rightarrow \infty} \left(\sum_{f=1}^p |x_{if} - x_{jf}|^h \right)^{1/h} = \max_f |x_{if} - x_{jf}|$$

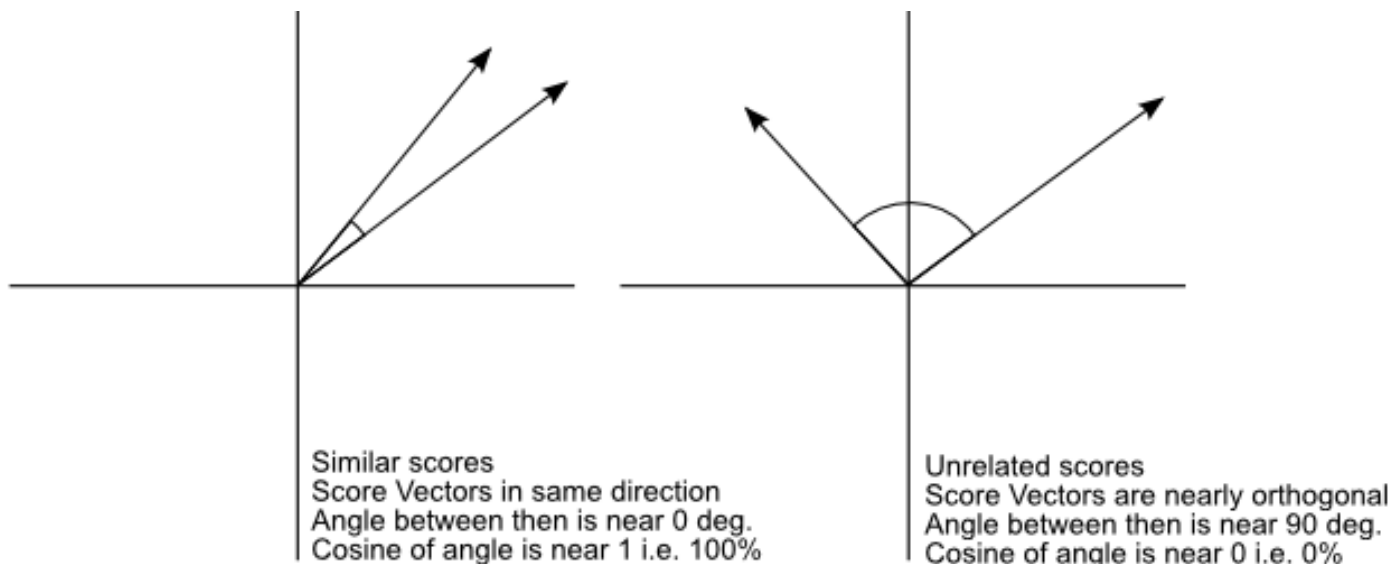
Cosine similarity

- A document can be represented by thousands of keywords in the document.

Document	team	coach	hockey	baseball	soccer	penalty	score	win	loss	season
<i>Document1</i>	5	0	3	0	2	0	0	2	0	0
<i>Document2</i>	3	0	2	0	1	1	0	1	0	1
<i>Document3</i>	0	7	0	$sim(d_1, d_2) = 0.94$			0	3	0	0
<i>Document4</i>	0	1	0				2	0	3	0

Cosine similarity

- Let d_1 and d_2 are two vectors (e.g., term-frequency vectors).
- Cosine similarity is **non-metric**:
$$\text{sim}(d_1, d_2) = \frac{d_1 \cdot d_2}{\|d_1\| \|d_2\|}$$
 - where \cdot is vector dot product, $\|d\|$ is the length of vector d
 - $\text{sim} = 0$ means no match, while $\text{sim} = 1$ means a complete match.



Measures for ordinal attributes

- The range of a numeric attribute can be mapped to an ordinal attribute f having M_f states.
 - E.g., temperate: cold ($-30^{\circ}\text{C} - 10^{\circ}\text{C}$), moderate ($-10^{\circ}\text{C} - 10^{\circ}\text{C}$), and warm ($10^{\circ}\text{C} - 30^{\circ}\text{C}$)
- Let M represent the number of possible ordered states, which define the ranking $1, \dots, M_f$
- Replace each x_{if} by its corresponding rank, $r_{if} \in \{1, \dots, M_f\}$
- Replace rank r_{if} of i^{th} object by $z_{if} = \frac{r_{if} - 1}{M_f - 1}$
- Continue with any measure for numeric attributes

Measures for ordinal attributes

Object Identifier	test-1 (nominal)	test-2 (ordinal)	test-3 (numeric)
1	code A	excellent	45
2	code B	fair	22
3	code C	good	64
4	code A	excellent	28

- test-2 = {fair, good, excellent}, i.e., $M_f = 3$
- The ranks of four objects are 3, 1, 2, and 3, respectively
- Map the rank 1 to 0.0, rank 2 to 0.5, and rank 3 to 1.0
- Dissimilarity matrix using Euclidean distance

$$\begin{bmatrix} 0 & & & \\ 1.0 & 0 & & \\ 0.5 & 0.5 & 0 & \\ 0 & 1.0 & 0.5 & 0 \end{bmatrix}$$

Measures for attributes of mixed types

- Suppose that the dataset has p attributes of mixed type.
- The distance between objects i and j is $d(i, j) = \frac{\sum_{f=1}^p \delta_{ij}^{(f)} d_{ij}^{(f)}}{\sum_{f=1}^p \delta_{ij}^{(f)}}$
 - $\delta_{ij}^{(f)} = 0$ if (1) x_{if} or x_{jf} is missing, or (2) $x_{if} = x_{jf} = 0$ and attribute f is asymmetric binary. Otherwise, $\delta_{ij}^{(f)} = 1$
 - If f is numeric: $d_{ij}^{(f)} = \frac{|x_{if} - x_{jf}|}{\max_h x_{hf} - \min_h x_{hf}}$, where h runs over all nonmissing objects for attribute f
 - If f is nominal or binary: $d_{ij}^{(f)} = 0$ if $x_{if} = x_{jf}$; otherwise, $d_{ij}^{(f)} = 1$
 - If f is ordinal: compute r_{if} and treat $z_{if} = \frac{r_{if} - 1}{M_f - 1}$ as numeric

Measures for attributes of mixed types

Dissimilarity matrix of test-1

$$\begin{bmatrix} 0 & & & \\ 1 & 0 & & \\ 1 & 1 & 0 & \\ 0 & 1 & 1 & 0 \end{bmatrix}$$

Dissimilarity matrix of test-2

$$\begin{bmatrix} 0 & & & \\ 1.0 & 0 & & \\ 0.5 & 0.5 & 0 & \\ 0 & 1.0 & 0.5 & 0 \end{bmatrix}$$

Dissimilarity matrix of test-3

$$\begin{bmatrix} 0 & & & \\ 0.55 & 0 & & \\ 0.45 & 1.00 & 0 & \\ 0.40 & 0.14 & 0.86 & 0 \end{bmatrix}$$

Object Identifier	test-1 (nominal)	test-2 (ordinal)	test-3 (numeric)
1	code A	excellent	45
2	code B	fair	22
3	code C	good	64
4	code A	excellent	28

- $\delta_{ij}^{(f)} = 1$ for each attribute f
- $d(3,1) = \frac{1(1)+1(0.50)+1(0.45)}{3} = 0.65$
- The resulting dissimilarity matrix

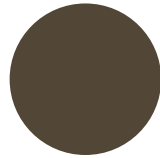
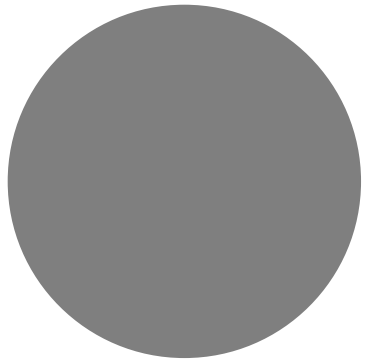
$$\begin{bmatrix} 0 & & & \\ 0.85 & 0 & & \\ 0.65 & 0.83 & 0 & \\ 0.13 & 0.71 & 0.79 & 0 \end{bmatrix}$$

Quiz 07: Jaccard coefficient

1. Calculate the similarity between these two observations, in which all the attributes are binary asymmetric.

IDs	fever	cough	breathing difficulty	fatigue	headache	loss of taste	sore throat
1	1	1	1	0	0	1	1
2	0	1	0	0	1	1	1

2. Explore the distance and similarity metrics supported in **scikit-learn**.
For each available metric, calculate the distance/similarity between the two observations above.



Correlation analysis



χ^2 statistics for correlation analysis

- Suppose attribute A has c distinct values and attribute B has r distinct values. There are n data tuples.
- Let (A_i, B_j) denote the joint event that $A = a_i$ and $B = b_j$.

- χ^2 statistic tests the null hypothesis, *A and B are independent*

$$\chi^2 = \sum_{i=1}^c \sum_{j=1}^r \frac{(o_{ij} - e_{ij})^2}{e_{ij}}$$

- o_{ij} : observed frequency (i.e., actual count) of $(A = a_i, B = b_j)$
- e_{ij} : expected frequency of (A_i, B_j) $e_{ij} = \frac{\text{count}(A=a_i) \times \text{count}(B=b_j)}{n}$
- The larger χ^2 value, the more likely the variables are related.

χ^2 statistics: An example

- Consider the below a contingency table.

	<i>male</i>	<i>female</i>	<i>Total</i>
<i>fiction</i>	250 (90)	200 (360)	450
<i>non_fiction</i>	50 (210)	1000 (840)	1050
Total	300	1200	1500

(Numbers in parenthesis are expected counts calculated based on the data distribution in the two categories)

- Are *gender* and *preferred_reading* correlated?

$$\chi^2 = \frac{(250 - 90)^2}{90} + \frac{(50 - 210)^2}{210} + \frac{(200 - 360)^2}{360} + \frac{(1000 - 840)^2}{840} = 507.93$$

- Two attributes are (strongly) correlated for the given group of people
- However, **correlation does not imply causality**.
 - # of hospitals and # of car-theft in a city are correlated
 - However, both are causally linked to the third variable – population.

χ^2 statistics: Contingency table

- A **contingency table** (or **crosstab**) displays the **frequency distribution** of two or more categorical variables.
- It helps to **analyze the relationship between variables**.

Values of the second variable

	male	female	Total
fiction	250 (90)	200 (360)	450
non_fiction	50 (210)	1000 (840)	1050
Total	300	1200	1500

Values of the first variable

Grand total

Marginal totals

χ^2 statistics: Contingency table

- Is a contingency table able to represent categorical variables that have more than two values?*

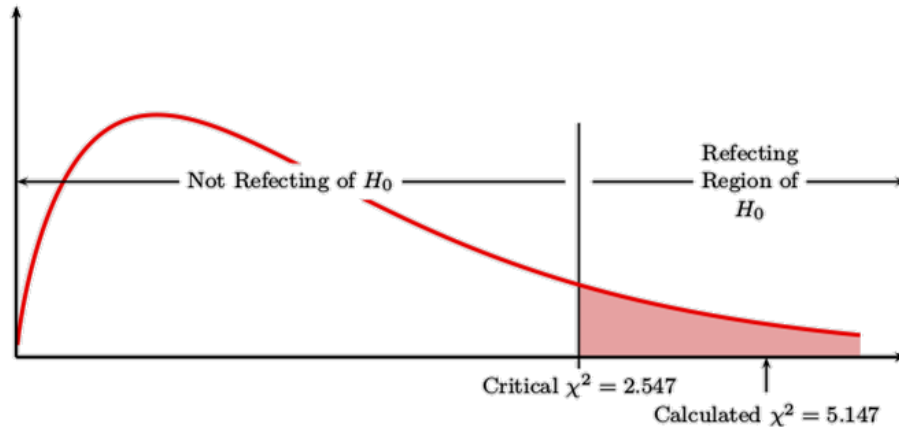
	18-29	30-49	50+	Row Total
Coffee	30	40	20	90
Tea	20	35	25	80
Juice	25	15	10	50
Column Total	75	90	55	220

χ^2 statistics: Contingency table

- Can a contingency table represent the relationship of more than two categorical variables?*

		Response (Y)	
Clinic (Z)	Drug Treatment (X)	Success	Failure
1	A	18	12
	B	12	8
2	A	2	8
	B	8	32

χ^2 statistics: An example



- The test is based on a significance level with a DOF of 1.
- If the hypothesis is denied, A and B are statistically correlated.

Degrees of Freedom	Probability										
	0.95	0.90	0.80	0.70	0.50	0.30	0.20	0.10	0.05	0.01	0.001
1	0.004	0.02	0.06	0.15	0.46	1.07	1.64	2.71	3.84	6.64	10.83
2	0.10	0.21	0.45	0.71	1.39	2.41	3.22	4.60	5.99	9.21	13.82
3	0.35	0.58	1.01	1.42	2.37	3.66	4.64	6.25	7.82	11.34	16.27
4	0.71	1.06	1.65	2.20	3.36	4.88	5.99	7.78	9.49	13.28	18.47
5	1.14	1.61	2.34	3.00	4.35	6.06	7.29	9.24	11.07	15.09	20.52
6	1.63	2.20	3.07	3.83	5.35	7.23	8.56	10.64	12.59	16.81	22.46
7	2.17	2.83	3.82	4.67	6.35	8.38	9.80	12.02	14.07	18.48	24.32
8	2.73	3.49	4.59	5.53	7.34	9.52	11.03	13.36	15.51	20.09	26.12
9	3.32	4.17	5.38	6.39	8.34	10.66	12.24	14.68	16.92	21.67	27.88
10	3.94	4.86	6.18	7.27	9.34	11.78	13.44	15.99	18.31	23.21	29.59
Non-significant									Significant		

χ^2 statistics: Degree of freedom (DOF)

- **DOF** is the number of independent values that can vary in a statistical calculation without breaking constraints.
- It can be calculated from the contingency table as follows.

$$DOF = (r - 1) \cdot (c - 1)$$

- r is the number of rows and c is the number of columns.
- E.g., in a table with 3 rows and 2 columns, $DOF = (3 - 1)(2 - 1) = 2$.

χ^2 statistics: Significance levels

- The **significance level α** is the probability threshold to decide **whether to reject the null hypothesis**.
- It represents the risk of rejecting a true null hypothesis.
- **Common significance levels**
 - **0.05 (5%)**: The most common choice, indicating a 5% risk of wrongly rejecting the null hypothesis.
 - **0.01 (1%)**: Used for stricter criteria, implying a 1% risk.
 - **0.10 (10%)**: Sometimes used in exploratory studies where a higher risk is acceptable.

Quiz 08: χ^2 statistics

1. Consider the data that relate the sex of children in families who have two children. Apply χ^2 statistics at the 0.001 significance level.

		First child		Total
		Male	Female	
Second child	Male	114 (120.54)	131 (124.46)	245
	Female	132 (125.46)	123 (129.54)	255
Total		246	254	500

Numbers out of parenthesis are actual counts from observations and numbers in parenthesis are expected counts calculated based on the data distribution in the two categories

2. How to perform χ^2 statistics using some **library in Python**?

Pearson correlation coefficient

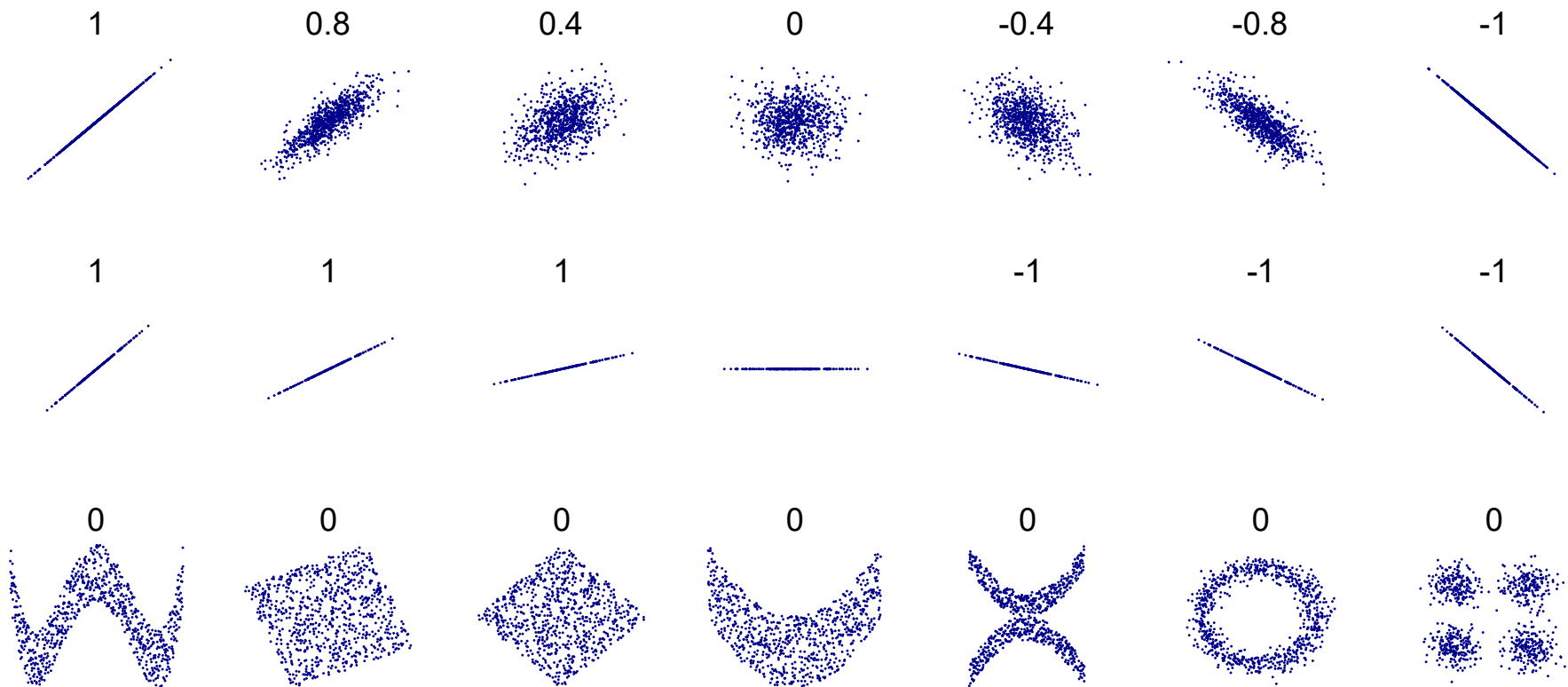
- Consider two numeric attributes A and B , and a set of n observations $\{(a_1, b_1), \dots, (a_n, b_n)\}$.
- **Pearson's product moment coefficient**

$$r_{A,B} = \frac{\sum_{i=1}^n (a_i - \bar{A})(b_i - \bar{B})}{n\sigma_A\sigma_B} = \frac{(\sum_{i=1}^n a_i b_i) - n\bar{A}\bar{B}}{n\sigma_A\sigma_B}$$

- $\bar{A}, \bar{B}, \sigma_A, \sigma_B$: means and standard deviations of A and B , respectively
- $\sum a_i b_i$: sum of the AB cross-product

$-1 \leftarrow r_{A,B}$	$r_{A,B} = 0$	$r_{A,B} \rightarrow 1$
Negative correlation	A and B are independent	Positive correlation

Pearson correlation coefficient



Several sets of (x, y) points, with the Pearson correlation coefficient of x and y for each set. The correlation reflects the noisiness and direction of a linear relationship (top row), but not the slope of that relationship (middle), nor many aspects of nonlinear relationships (bottom). N.B.: the figure in the center has a slope of 0 but in that case the correlation coefficient is undefined because the variance of Y is zero. ([Wikipedia](https://en.wikipedia.org/wiki/Pearson_product-moment_correlation_coefficient))

Covariance analysis

- The **covariance between A and B** is defined as

$$Cov(A, B) = \frac{\sum_{i=1}^n (a_i - \bar{A})(b_i - \bar{B})}{n} = E(A \cdot B) - \bar{A}\bar{B}$$

- where $E(A) = \bar{A} = \frac{\sum_{i=1}^n a_i}{n}$ and $E(B) = \bar{B} = \frac{\sum_{i=1}^n b_i}{n}$ are the expected values of A and B

$Cov(A, B) > 0$	$Cov(A, B) < 0$	$Cov(A, B) = 0$
Positive covariance	Negative covariance	A and B are independent

- Covariance vs. correlation: $r_{A,B} = \frac{Cov(A,B)}{\sigma_A \sigma_B}$

Covariance analysis: An example

- If the stocks are affected by the same industry trends, will their prices rise or fall together?

Stock Prices for *AllElectronics* and *HighTech*

<i>Time point</i>	<i>AllElectronics</i>	<i>HighTech</i>
t1	6	20
t2	5	10
t3	4	14
t4	3	5
t5	2	5

- $E(\textit{AllElectronics}) = \frac{6+5+4+3+2}{5} = \frac{20}{5} = \4
- $E(\textit{HighTech}) = \frac{20+10+14+5+5}{5} = \frac{54}{5} = \10.80
- $\textit{Cov}(\textit{AllElectronics}, \textit{HighTech}) = \frac{6 \times 20 + 5 \times 10 + 4 \times 14 + 3 \times 5 + 2 \times 5}{5} - 4 \times 10.80 = 7$
- Therefore, a positive covariance indicates that stock prices for both companies rise together

Quiz 09: Correlation metrics

1. Consider the following data table, in which there are five tuples of two attributes, A and B.

Calculate the Pearson correlation coefficient and Covariance between A and B.

No.	Attributes	
	A	B
1	19	16
2	25	10
3	13	26
4	12	29
5	16	20

2. How to compute the above metrics using **pandas**? Show the results for the above data.

References

- Jiawei Han, Micheline Kamber, and Jian Pei, 2011. Data Mining: Concepts and Techniques (3rd ed.). Morgan Kaufmann Publishers Inc. Chapter 2 and Chapter 3.

...the end.

