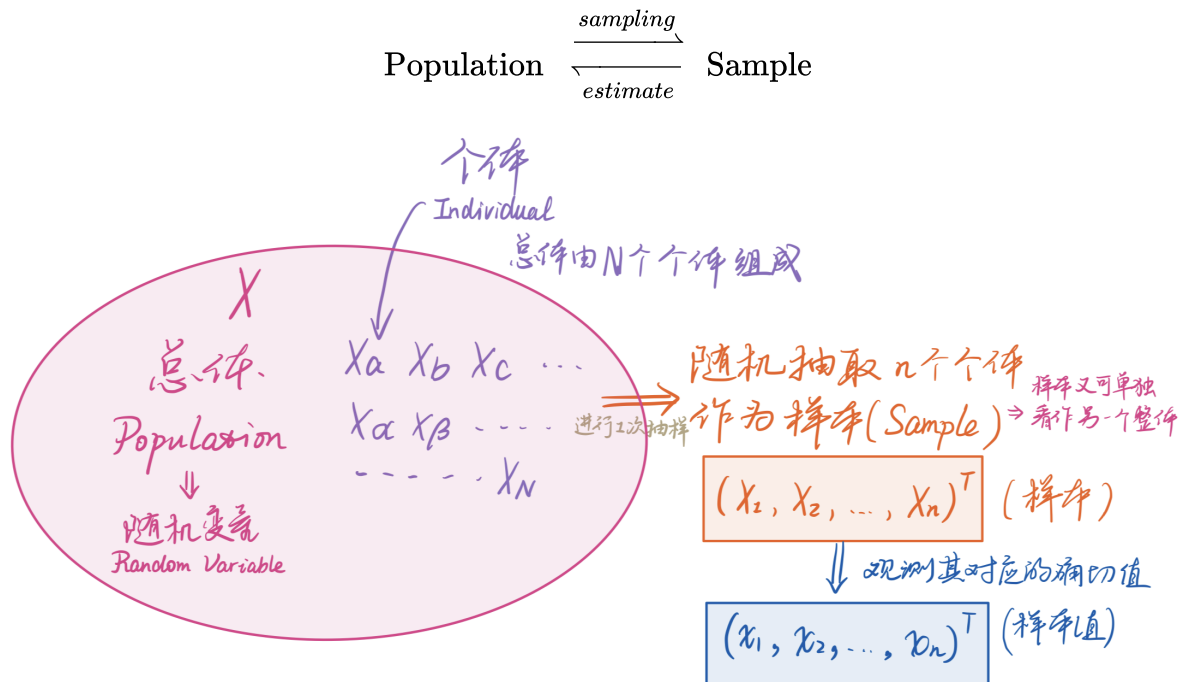


Python Statistical Analysis

There are some fundamental symbol and concept need to know first.



Statistic (统计量)	Population (总体) - N 个	Sample (样本) - n 个
Mean	μ	$\bar{X} = \hat{\mu}$
Variance	σ^2	$s^2 = \hat{\sigma}^2$
Proportion	π	p
Correlative Coefficeince	ρ	r

```

1  import numpy as np
2
3  # Convert multiple-rows data to the MATRIX of COLUMN VECTORS
4  np.matrix([<data_1>, <data_2>, ...]).T
5
6  > input:  data_1 = [0, 1, 2]
7            data_2 = [9, 8, 7]
8            print(np.matrix([data_1, data_2])) # Left
9            print(np.matrix([data_1, data_2]).T) # Right
10
11  > output: +---+---+---+      +---+---+
12            | 0 | 1 | 2 |      | 0 | 9 |
13            +---+---+---+      +---+---+
14            | 9 | 8 | 7 |      | 1 | 8 |
  
```

```

14      +---+---+---+      +---+---+
15      | 2 | 7 |
16      +---+---+
17
18
19 # Ravel the 2-D data to 1-D data
20 <data>.ravel
21
22 > input:  m = np.array([[1, 2, 3],
23                        [0, 1, 2]])
24           print(m.mean(axis=0))
25           print(m.mean(axis=0).ravel())
26 > output: [[0.5, 1.5, 2.5]]
27           [0.5, 1.5, 2.5]

```

1. Data Descriptive Analysis · 描述性统计分析

1.1. Descriptive Statistics · 描述性统计量

1.1.1. Measures of Location and Dispersion · 位置和分散程度的度量

There are 8 important statistics used to measure the location and dispersion, which are **Mean**, **Median**, **Percentile**, **Variation**, **Standard Deviation**, **Standard Error of Mean**, **Coefficient of Variation**, and **Range**.

1. Mean (均值, $\mathbb{E}[X] = \bar{X} = \langle X \rangle$)

The average of all entries.

$$\mu = \frac{1}{N} \sum_{i=1}^N X_i$$

$$\hat{\mu} = \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \quad (1.1)$$

e.g.

$$X = \{1, 2, 3, 4, 5\}$$

$$\bar{X} = \frac{1 + 2 + 3 + 4 + 5}{5} = \frac{15}{5} = 3$$

```

1  import script.stats as st
2
3  # 1. Use numpy
4  #   axis = 0   → column-wise
5  #           = 1   → row-wise
6  #           = none → all entries
7  np.mean(<data>, axis=)
8  <ndarray-like data>.mean()
9
10 # 2. Use scipy.stats
11 #   Compute the mean of the closed interval
12 st.tmean(<data>, <interval as (a, b)>)

```

2. Median (中位数)

The middle number of all entries. **Data must be sorted in ascending order** (no need to sort first with Python).

In **Symmetric Distribution** (like t -distribution and normal-distribution, 对称分布), the median is very close to the mean; while in **Skewed Distribution** (like F -distribution, 偏态分布), the difference between median and mean is relatively large.

$$m_e = \begin{cases} x_{\frac{N+1}{2}} & , \text{if } N \text{ is odd} \\ \frac{1}{2} \left(x_{\frac{N}{2}} + x_{\frac{N+1}{2}} \right) & , \text{if } N \text{ is even} \end{cases} \quad (1.2)$$

e.g.

$$X = \{1, 2, 3, 4, 5, 6\}$$

$$m_e = \frac{1}{2} (x_3 + x_4) = \frac{3 + 4}{2} = 3.5$$

```
1 # reverse: False by default (descending)
2 sorted(<data>, reverse=True) # Order data by ascending
3
4 np.median(<data>)
```

3. Percentile (百分位数)

Represent specific-location data. **Data must be sorted in ascending order** (no need to sort first with Python). The symbol $\lceil \cdot \rceil$ means rounded up (e.g., $\lceil 2.1 \rceil = \lceil 2.9 \rceil = 3$).

$$m_p = \begin{cases} x_{\lceil N \times \text{perc} \rceil} & , \text{if } N \times \text{perc} \notin \mathbb{N}^+ \\ \frac{1}{2} (x_{N \times \text{perc}} + x_{N \times \text{perc} + 1}) & , \text{if } N \times \text{perc} \in \mathbb{N}^+ \end{cases} \quad (1.3)$$

e.g.1

$$X = \{1, 2, 3, 4, 5\}$$

$$m_{50\%} = x_{\lceil 5 \times 50\% \rceil} = x_{\lceil 2.5 \rceil} = x_3 = 3$$

e.g.2

$$X = \{1, 2, 3, 4, 5, 6\}$$

$$m_{50\%} = \frac{(x_{50\% \times 6} + x_{50\% \times 6 + 1})}{2} = \frac{(x_3 + x_4)}{2} = \frac{3 + 4}{2} = 3.5$$

```
1 np.quantile(<data>, <percentile>)
```

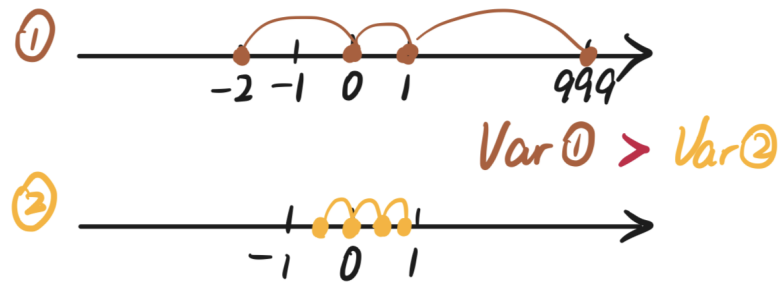
4. Variance (Var, 方差, $\text{Var}[X]$)

Measure the dispersion of data, namely the degree of the samples dispersion from the mean. **The more concentrated the data, the smaller the variance, vice versa.**

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (X_i - \mu)^2$$

$$\begin{cases} \hat{\sigma}_{biased}^2 = s_{biased}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 \\ \hat{\sigma}_{unbiased}^2 = s_{unbiased}^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 \end{cases} \quad (1.4)$$

e.g., The following picture can easily figure the relationship between *variance* and *dispersion* of data.



```
1 # Biased Estimator of variance
2 np.var(<data>)
3
4 # Unbiased Estimator of variance
5 st.tvar(<data>)
```

There are two kinds of estimation: Unbiased Estimation (无偏估计) and Biased Estimation (有偏估计).

If we want to discern an unknown distribution (i.e., μ, σ^2 are all unknown), we need to take some samples to estimate its $\{\mu, \sigma^2\}$. Mean can be easily compute by (1.1), but variance is a little complex for which refers the conception of the **Freedom Degree** (自由度).

Consider an extreme case that we **take only one sample** $X = \{X_1\}$, so we can easily estimate that

$$\text{the population mean estimator } \hat{\mu} = \text{the sample mean } \bar{X} = \frac{X_1}{1}$$

Estimator: an estimation **function** of samples (i.e., \bar{X}).

Estimate: a **certain value** compute by substitute one sample to the estimator (i.e., $\frac{1+2}{2} = 1.5$).

But when we want to estimate the population variance σ^2 by the sample variance $s^2 = \hat{\sigma}^2$ with $\hat{\mu}$, we will find something counterintuitive:

- if we compute $\hat{\sigma}^2$ by divide n intuitively, we will get $\hat{\sigma}^2 = \frac{(X_1 - X_1)^2}{1} = 0$ which is counterintuitively, because the only one element does not exist $\hat{\sigma}^2$, so we cannot compute $\hat{\sigma}^2$ by only one element;

- while if we compute $\hat{\sigma}^2$ by divide $(n - 1)$ counterintuitively, we will get $\hat{\sigma}^2 = \frac{(X_1 - X_1)^2}{1 - 1} = \frac{0}{0}$ which is cannot be computed but intuitively. This means that although the numerator looks like having an item $(X_1 - X_1)^2$, **the information it contains is actually 0 (i.e., freedom degree $df = 0$).**

Expand this case to two samples: when we have two samples $X = \{X_1, X_2\}$, we can get $\hat{\mu} = \bar{X} = \frac{X_1 + X_2}{2}$, $\hat{\sigma}^2 = s^2 = \frac{(X_1 - \bar{X})^2 + (X_2 - \bar{X})^2}{2}$. However actually, we can find that $(X_1 - \bar{X}) = -(X_2 - \bar{X})$, namely one of them **is not free**, so the information it contains is actually 1 (i.e., freedom degree $df = 1$). **By analogy, when we have n samples, we get $df = n - 1$** , which means we only calculated $(n - 1)$ deviances from the population, so that:

- To **estimate the population variance σ^2 closest**, we **calculate the sample variance s^2 by divide $(n - 1)$** , which named **unbiased estimator of the population variance $\hat{\sigma}^2$** , where $\mathbb{E}(s^2) = \sigma^2$;
- And if we calculate s^2 by n , we will get **biased estimator of the population variance $\hat{\sigma}^2$** , where $\mathbb{E}(s^2) \neq \sigma^2$.

5. Standard Deviation (SD, 标准差)

The square root of variation.

$$\sigma = \sqrt{\frac{1}{N} \sum_{i=1}^N (X_i - \mu)^2}$$

$$\begin{cases} \hat{\sigma}_{biased} = s_{biased} = \sqrt{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2} \\ \hat{\sigma}_{unbiased} = s_{unbiased} = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2} \end{cases} \quad (1.5)$$

```
1 # Biased Estimator of standard deviation
2 np.std(<data>)
3
4 # Unbiased Estimator of standard deviation
5 st.tstd(<data>)
```

6. Standard Error of the Mean (SEM, 标准误)

Measure the dispersion of the sample mean for every sample, and the deviation of the population mean μ and the sample mean \bar{X} . **The smaller the SEM, the smaller the deviation between μ and \bar{X} , vice versa.**

$$\hat{\sigma}_{SEM} = \sqrt{\frac{1}{m-1} \sum_{j=1}^m (\bar{X}_j - \bar{\bar{X}})^2} \cong \frac{s_{j_unbiased}}{\sqrt{n}} \quad (1.6)$$

where $i = \{1, \dots, n\}$ means the number of entries for every sample, $j = \{1, \dots, m\}$ means the number of samples. In general, we sample one time and use *this sample SD* divided by \sqrt{n} , as shown on the right of (1.6).

e.g., Suppose we sample three times, every sample have four entries. The data are shown in the table below.

sample j entry i	A	B	C
1	0	0	0
2	1	2	4
3	2	4	8
4	3	6	12
\bar{X}_j	1.5	3	6
$s_{j_unbiased}$	1.29	2.58	5.16

$$\bar{\bar{X}} = \frac{1.5 + 3 + 6}{3} = 3.5$$

$$\hat{\sigma}_m = \sqrt{\frac{(1.5 - 3.5)^2 + (3 - 3.5)^2 + (6 - 3.5)^2}{2}} \approx 2.29 \cong \frac{1.29}{\sqrt{4}} = 0.645$$

```
1 # Calculate the SEM by hand
2 sem = st.tstd(<data>) / np.sqrt(<the number of entries for this
  sample>)
```

7. Coefficient of Variation (CV, 变异系数)

Measure the dispersion of data by eliminate the effects of *dimension* (尺度) and *scale* (量纲).

$$CV = \frac{\sigma}{\mu} \quad (1.7)$$

```
1 # Calculate the CV by hand
2 cv = st.tstd(<data>) / np.mean(<data>)
```

8. Range (极差, 全距)

The longest length of data.

$$\text{Range} = \max - \min \quad (1.8)$$

```
1 # Calculate the Range by hand
2 Range = np.max(<data>) - np.min(<data>)
```

1.1.2. Measures of Relationships · 关系度量

1. Variance-Covariance Matrix (方差-协方差矩阵) & Correlation Coefficient Matrix (相关系数矩阵)

Covariance: Measure the changing trends between the 2 or more variables. This matrix is symmetric along the diagonal.

$$\begin{cases} \hat{\sigma}_{xx}^2 = s_{xx}^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 = \text{Var}(X) \\ \hat{\sigma}_{yy}^2 = s_{yy}^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})^2 = \text{Var}(Y) \\ \hat{\sigma}_{xy}^2 = s_{xy}^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}) = s_{yx}^2 = \hat{\sigma}_{yx}^2 \end{cases} \quad (1.9)$$

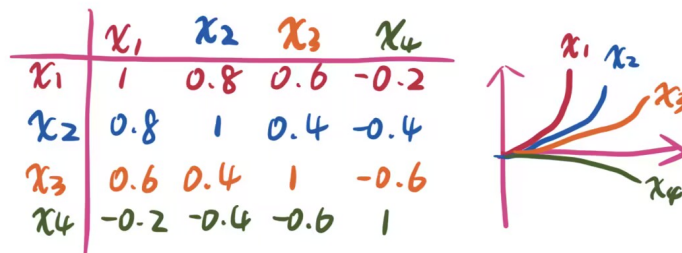
Correlation Coefficient: **Normalize (归一化) the covariance** to measure the degree of linear correlation between 2 or more variables.

$$\rho_{XY} = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} \in [-1, 1] \quad (1.10)$$

The meaning of ρ is as follows



and the visualization is as follows



```
1 # Variance-Covariance Matrix
2 np.cov(<data>) # The data form must be a ROW VECTOR
3
4
5 # Correlation Coefficient Matrix
6 # 1. Use Numpy
7 np.corrcoef(<data>) # The data form must be a ROW VECTOR
8
9 # 2. Use Pandas
10 pd.DataFrame.corr(<data>) # The data form must be a ROW VECTOR
```

1.1.3. Measures of Distribution Shapes · 分布形状的度量

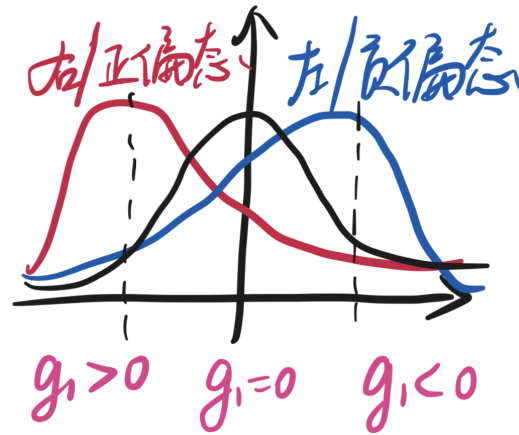
1. Skewness (偏度)

Measure the degree of asymmetry in the data. The more the SK is to 0, the more the curve fits the normal distribution.

$$g_1 = SK(X) = \frac{n \sum_{i=1}^n (X_i - \bar{X})^3}{(n-1)(n-2)s^3} = \frac{n^2 \mu_3}{(n-1)(n-2)s^3}$$

where s means the Standard Deviation, μ_3 means Third Order Central Moment (三阶中心矩), $\mu_3 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^3$.

The distribution shapes under different skewness are as follows



```
1 # 1. Use Pandas, calculate the modified skewness (Unbiased Estimation)
2 data_s = pd.Series(<data>)
3 data_s.skew()
4
5 # 2. Use Scipy
6 # bias = False → Modified (Unbiased)
7 #      = True  → Unmodified (Biased) (Default)
8 st.skew(<data>, bias=)
```

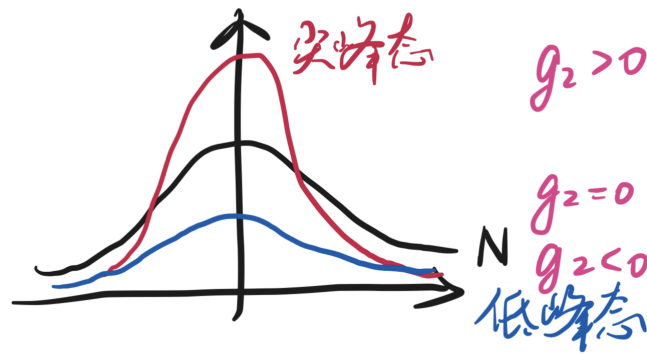
2. Kurtosis (峰度)

Measure the sharpness of the Probability Density Function (PDF, 概率密度函数) at μ .

$$\begin{aligned} g_2 = K(X) &= \frac{n(n+1) \sum_{i=1}^n (X_i - \bar{X})^4}{(n-1)(n-2)(n-3)s^4} - 3 \frac{(n-1)^2}{(n-2)(n-3)} \\ &= \frac{n^2(n+1)\mu_4}{(n-1)(n-2)(n-3)s^4} - 3 \frac{(n-1)^2}{(n-2)(n-3)} \end{aligned}$$

where s means the Standard Deviation, μ_4 means Fourth Order Central Moment (四阶中心矩), $\mu_4 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^4$.

The distribution shapes under different kurtosis are as follows



```

1 # 1. Use Pandas, calculate the modified kurtosis (Unbiased Estimation)
2 data_s = pd.Series(<data>)
3 data_s.kurt()
4
5 # 2. Use Scipy
6 # bias = False → Modified (Unbiased)
7 #      = True  → Unmodified (Biased) (Default)
8 st.kurtosis(<data>, bias=)

```

1.1.4. Summary of Data Characteristics · 数据特性的总括

Data characteristics usually means the **Global Extremum**, **Mean**, **Unbiased Variance**, **Modified Skewness**, **Modified Kurtosis**, and **Distribution**. The *descriptive statistics* can be calculated by the above content, and the *distribution* can be tested by **Normality Test** (like **Shapiro Test**) or **Distribution Fit Test** (like **Kolmogorov-Smirnov Test**)

```

1 # Check the decriptive statistics
2 st.describe(<data>, bias=False)
3
4 # Test the distribution
5 # 1. Normality Test
6 st.shapiro(<data>)
7 # 2. Distribution Fit Test
8 st.kstest(<data>, <name of the dist as "t">, <df of the dist as (3,)>)
9
10 # p.s. Generate an F-distribution with df=(2, 9)
11 st.f.rvs(size=, dfn=2, dfd=9)

```

1.2. Data Distribution · 数据分布

1.2.1. Fundamental Concepts · 基础概念

1. **Probability Space** ((Ω, \mathcal{F}, P) , 概率空间),
Sample Space (Ω , 样本空间),
Random Events ($A \in \mathcal{F}$, 随机事件),
Random Variable (RV, 随机变量)

asd

- 2. **Population** (总体),
Individual (个体),
Sample (样本),
Parameter Space (参数空间),
Distribution Family (分布族),
Statistics (统计量),
Sampling Distribution (抽样分布)

asd

- 3. **Probability Distribution Function** (Culmulative Distribution Function, CDF, 累积分布函数)

asd

- 4. **Discrete Random Variable** (离散型随机变量)

asd

- 5. **Continuous Random Variable** (连续型随机变量)

asd

- 6. **Mathematic Expectation** ($\mathbb{E}(\cdot)$, 数学期望)

asd

- 7. **Moment** (矩, 可理解为一种距离)

asd

1.2.2. Discrete Probability Distribution · 离散型概率分布

- 1. **Binomial Distribution** (二项分布)

asd

- 2. **Poisson Distribution** (泊松分布)

asd

1.2.3. Continuous Probability Distribution · 连续型概率分布

- 1. **Normal Distribution** (正态分布)

asd

- 2. **t-distribution** (Student's t distribution, t 分布)

asd

- 3. **Gamma Distribution** (伽马分布)

asd

- 4. **Chi-Square Distribution** (卡方分布)

asd

- 5. **F-distribution** (F 分布)

asd

1.3. Histogram, Experience Distribution and QQ Graph · 直方图, 经验分布函数与 QQ 图

1.3.1. Histogram and Kernel Density Estimation · 直方图与核密度估计

1.3.2. Experience Distribution · 经验分布函数

1.3.3. QQ Graph and Stem-and-Leaf Display · QQ 图与茎叶图

1.4. Multivariate Data · 多元数据

1.4.1. Numerical Features of Multivariate Data · 多元数据的数字特征

asd

1.4.2. Graphical Representation of Multivariate Data · 多元数据的图形表示

asd