

法律声明

■ 本课件包括演示文稿、示例、代码、题库、视频和声音等内容，北风网和讲师拥有完全知识产权；只限于善意学习者在本课程使用，不得在课程范围外向任何第三方散播。任何其他人或者机构不得盗版、复制、仿造其中的创意和内容，我们保留一切通过法律手段追究违反者的权利。

■ 课程详情请咨询

◆ 微信公众号：北风教育

◆ 官方网址：<http://www.ibeifeng.com/>



人工智能之机器学习

数学基础、Python科学计算库回顾

主讲人：Gerry

上海育创网络科技有限公司



课程要求

■ 课上课下 “九字” 真言

- ◆ 认真听，善摘录，勤思考
- ◆ **多温故，乐实践**，再发散

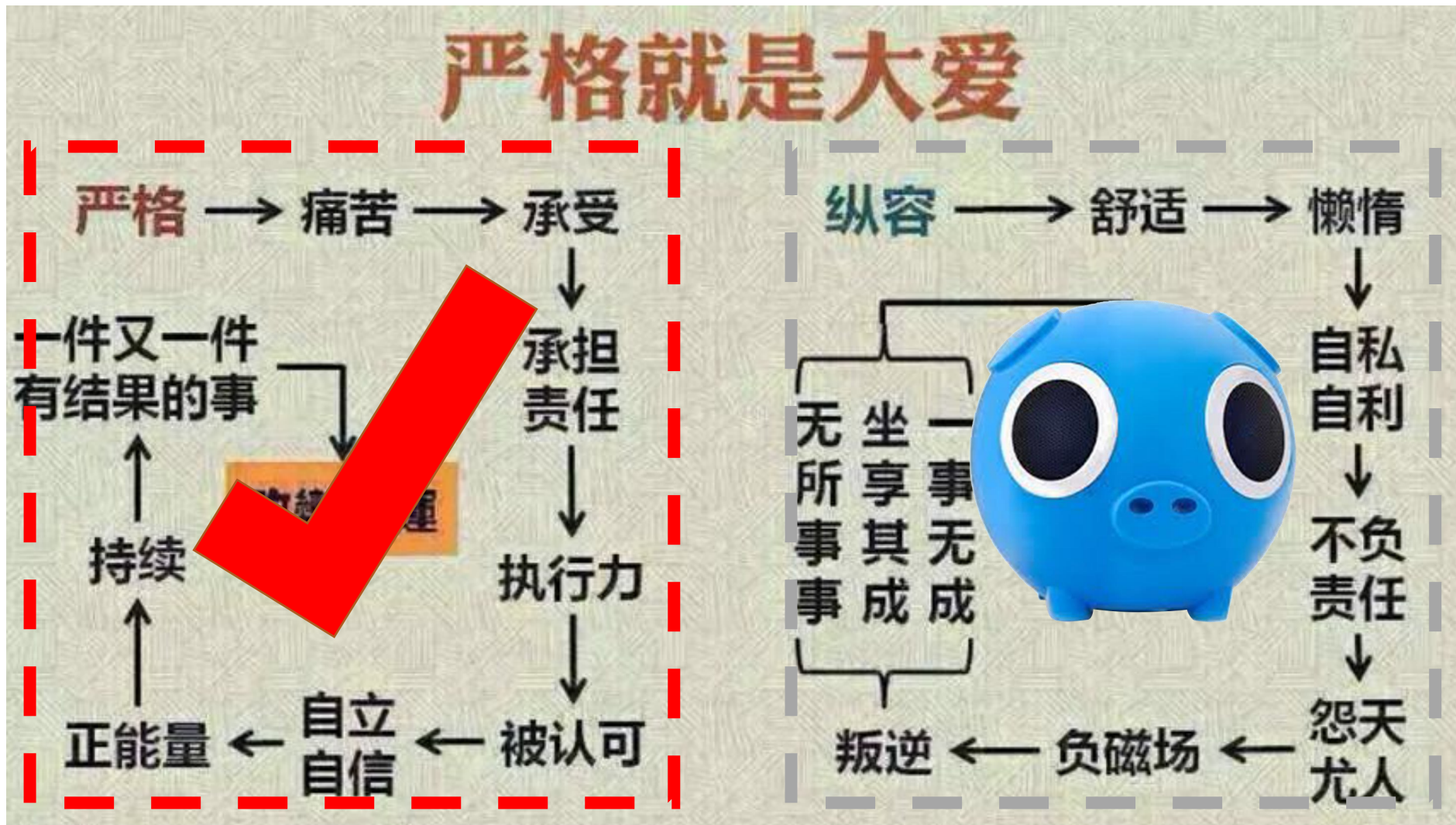
■ 四不原则

- ◆ **不懒散惰性，不迟到早退**
- ◆ **不请假旷课，不拖延作业**

■ 一点注意事项

- ◆ 违反 “四不原则”，不包就业和推荐就业

严格是大爱



寄语



做别人不愿做的事，
做别人不敢做的事，
做别人做不到的事。

课程内容

- 数学基础回顾
- 梯度下降法讲解
- Python科学计算库回顾

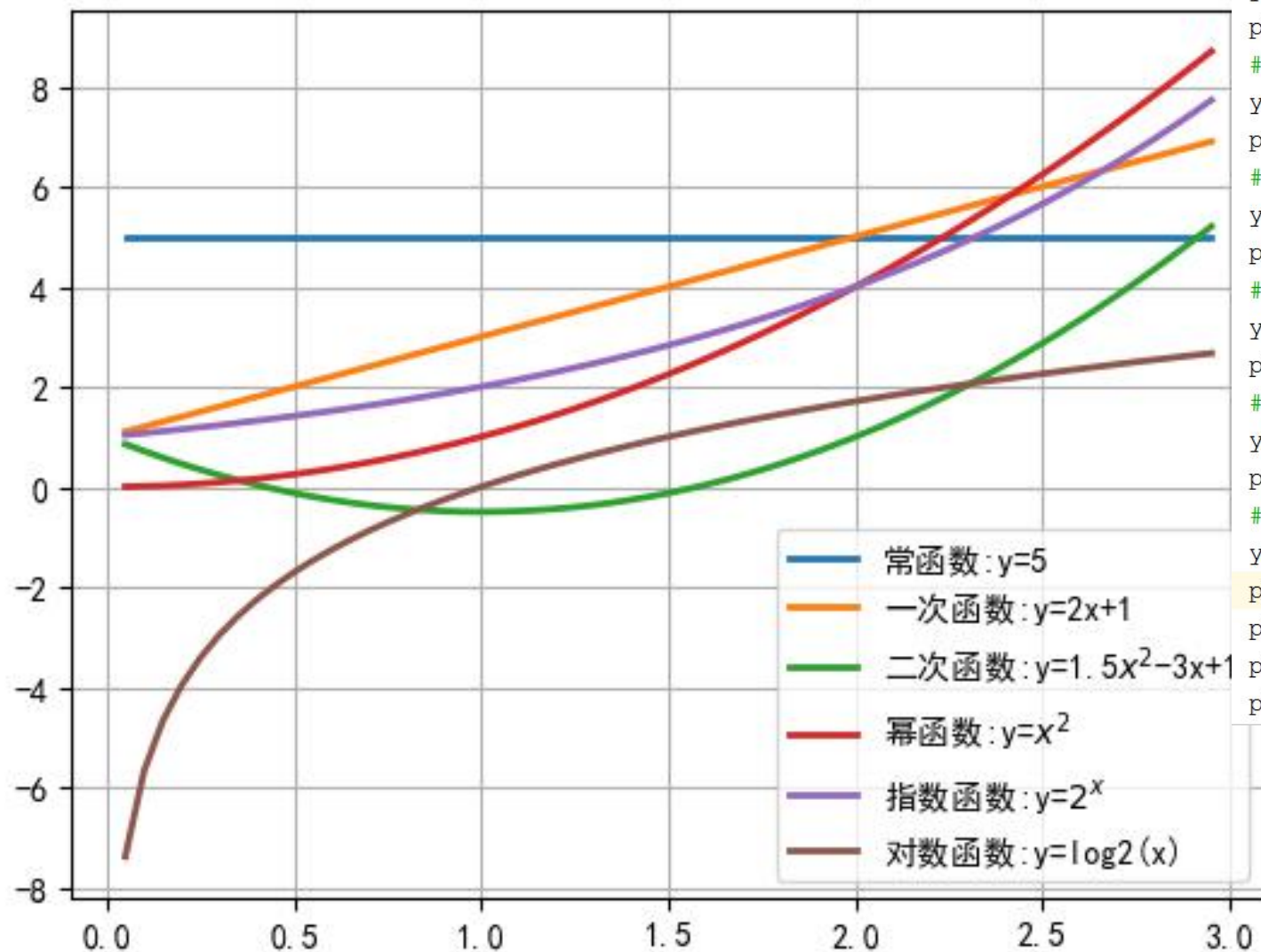
数学知识回顾

- 常见函数
- 导数、梯度
- Taylor公式
- 联合概率、条件概率、全概率公式、贝叶斯公式
- 期望、方差、协方差
- 大数定理、中心极限定理
- 最大似然估计(MLE)
- 向量、矩阵的运算
- 向量、矩阵的求导
- SVD、QR分解

常见函数

- 常函数: $y = C$ 一次函数: $y = ax + b$
- 二次函数: $y = ax^2 + bx + c$ 幂函数: $y = x^a$
- 指数函数: $y = a^x$, a 的取值范围为: $a > 0 \& a \neq 1$
- 对数函数: $y = \log_a(x)$, a 的取值范围为: $a > 0 \& a \neq 1$

常见函数



```
x = np.arange(0.05, 3, 0.05)
# 常函数
y1 = [5 for i in x]
plt.plot(x, y1, linewidth=2, label='常函数:y=5')
# 一次函数
y2 = [2 * i + 1 for i in x]
plt.plot(x, y2, linewidth=2, label='一次函数:y=2x+1')
# 二次函数
y3 = [1.5 * i * i - 3 * i + 1 for i in x]
plt.plot(x, y3, linewidth=2, label='二次函数:y=1.5x^2-3x+1')
# 幂函数
y4 = [math.pow(i, 2) for i in x]
plt.plot(x, y4, linewidth=2, label='幂函数:y=x^2')
# 指数函数
y5 = [math.pow(2, i) for i in x]
plt.plot(x, y5, linewidth=2, label='指数函数:y=2^x')
# 对数函数
y6 = [math.log(i, 1.5) for i in x]
plt.plot(x, y6, linewidth=2, label='对数函数:y=log2(x)')
plt.legend(loc='lower right')
plt.grid(True)
plt.show()
```

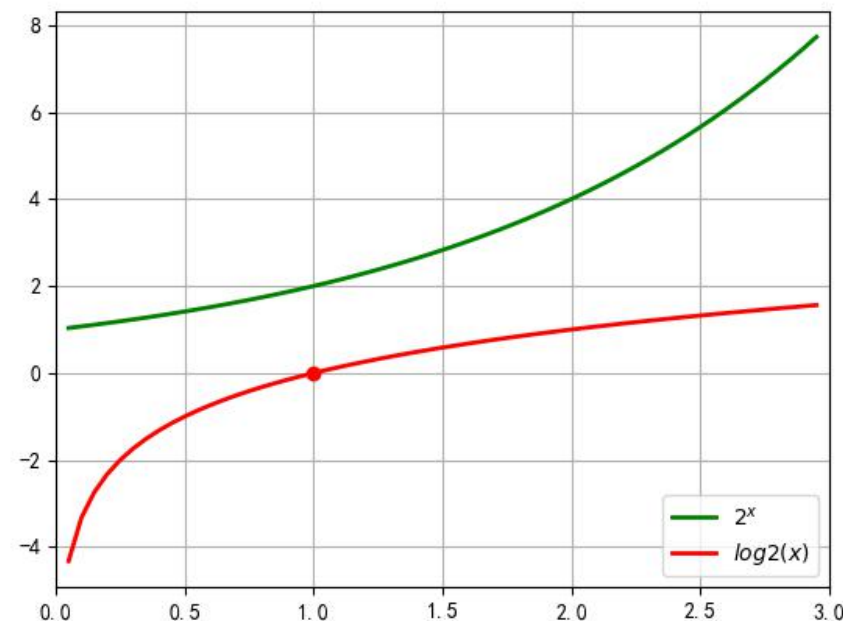
对数函数、指数函数

- 对数函数： $y = \log_a(x), a > 0 \text{ 且 } a \neq 1$
- 指数函数： $y = a^x, a > 0 \text{ 且 } a \neq 1$

$$\log_a x + \log_a y = \log_a (x * y)$$

$$\log_a x - \log_a y = \log_a \frac{x}{y}$$

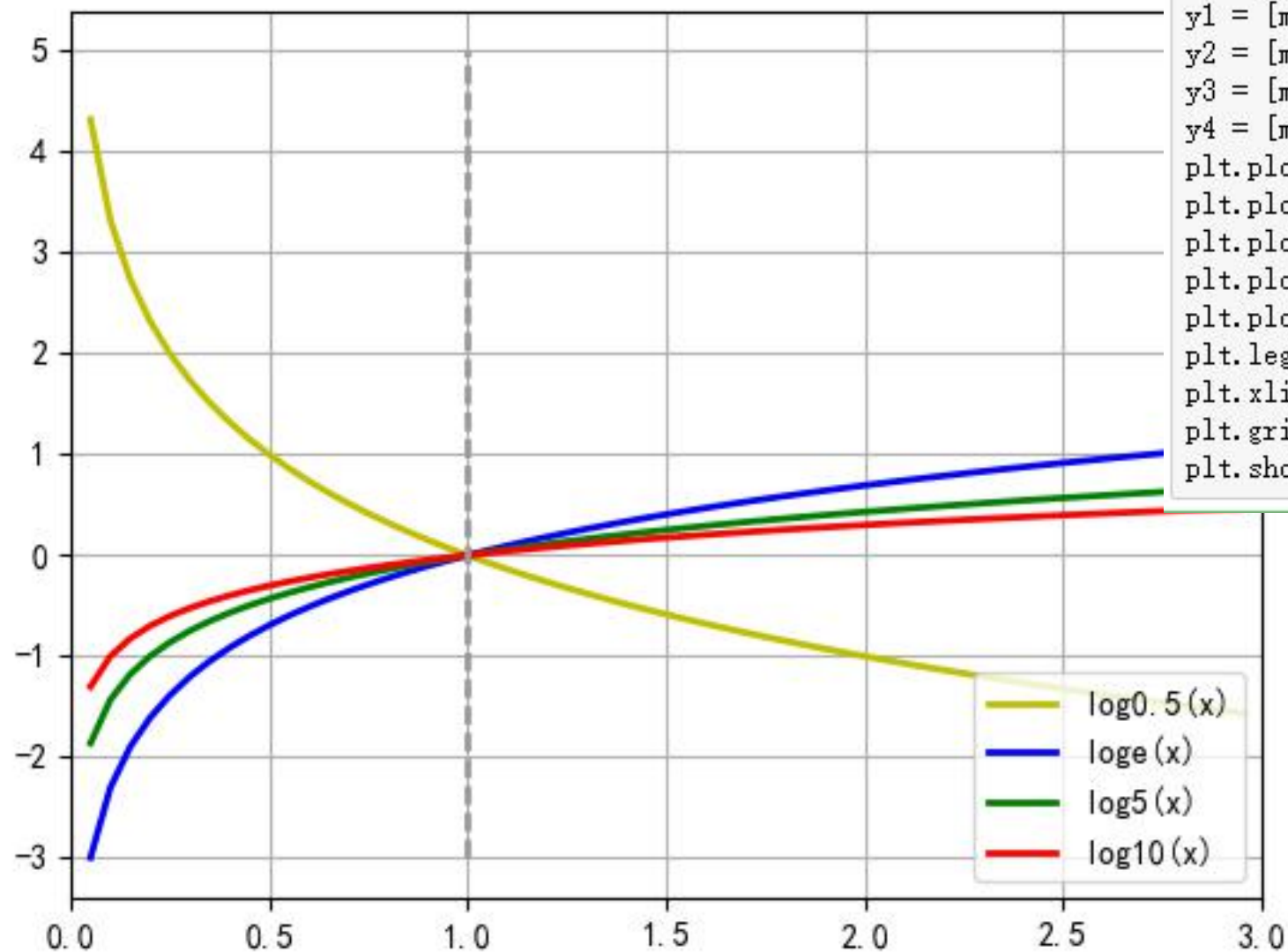
$$\frac{a^x}{a^y} = a^{x-y} \quad a^x \cdot a^y = a^{x+y}$$



$$\log_a(1) = 0$$

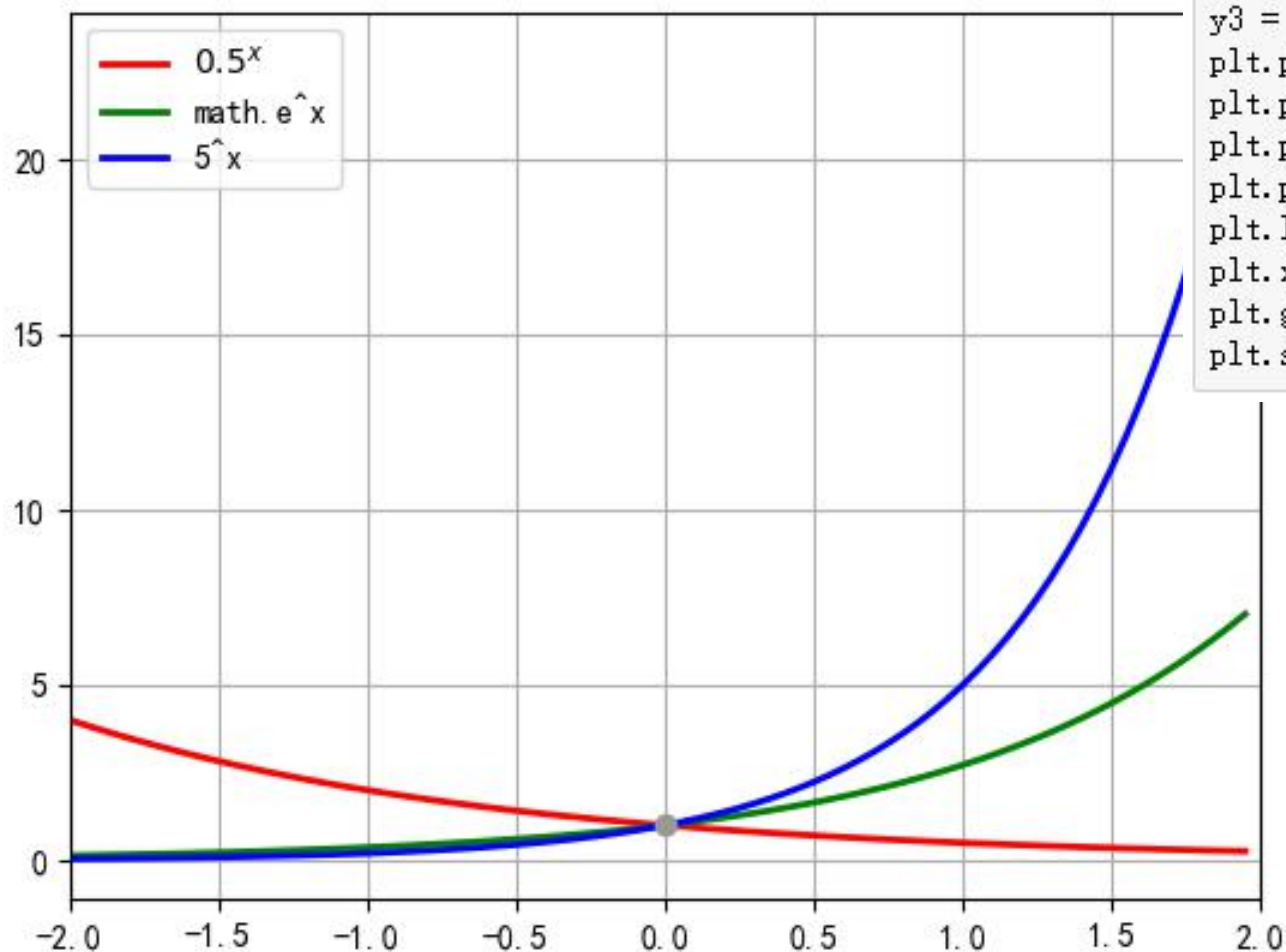
$$a^0 = 1$$

对数函数、指数函数



```
x = np.arange(0.05, 3, 0.05)
y1 = [math.log(i, 0.5) for i in x]
y2 = [math.log(i, math.e) for i in x]
y3 = [math.log(i, 5) for i in x]
y4 = [math.log(i, 10) for i in x]
plt.plot(x, y1, linewidth=2, color='y', label='log0.5(x)')
plt.plot(x, y2, linewidth=2, color='b', label='loge(x)')
plt.plot(x, y3, linewidth=2, color='g', label='log5(x)')
plt.plot(x, y4, linewidth=2, color='r', label='log10(x)')
plt.plot([1,1], [-3, 5], '-', color='#999999', linewidth=2)
plt.legend(loc='lower right')
plt.xlim(0, 3)
plt.grid(True)
plt.show()
```

对数函数、指数函数

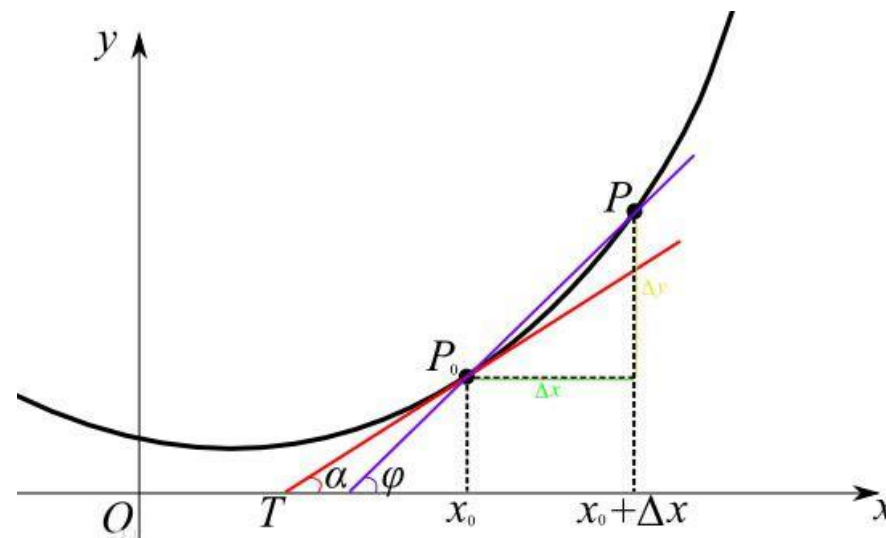


```
x = np.arange(-2, 2, 0.05)
y1 = [math.pow(0.5, i) for i in x]
y2 = [math.pow(math.e, i) for i in x]
y3 = [math.pow(5, i) for i in x]
plt.plot(x, y1, linewidth=2, color='r', label='$0.5^x$')
plt.plot(x, y2, linewidth=2, color='g', label='math.e^x')
plt.plot(x, y3, linewidth=2, color='b', label='5^x')
plt.plot([0], [1], 'o', color='#999999', linewidth=2)
plt.legend(loc='upper left')
plt.xlim(-2, 2)
plt.grid(True)
plt.show()
```

导数

- 一个函数在某一点的导数描述了这个函数在这一点附近的**变化率**，也可以认为是函数在某一点的导数就是该函数所代表的曲线在这一点**的切线斜率**。导数值越大，表示函数在该点处的变化越大。
- 定义：当函数 $y=f(x)$ 在自变量 $x=x_0$ 上产生一个增量 Δx 时，函数输出值的增量 Δy 和自变量增量 Δx 之间的比值在 Δx 趋近与0的时候存在极限值 a ，那么 a 即为函数在 x_0 处的导数值。

$$\lim_{\Delta x \rightarrow 0} \frac{f(x_0 + \Delta x) - f(x_0)}{\Delta x} = \lim_{\Delta x \rightarrow 0} \frac{\Delta y}{\Delta x}$$



导数

- **导数**就是曲线的斜率，是曲线变化快慢的一个反应。
- **二阶导数**是斜率变化的反应，表现曲线是**凹凸性**。

$$y = f(x)$$

$$y' = f'(x) = \frac{\partial f(x)}{\partial x} = \frac{dy}{dx} = \lim_{\Delta x \rightarrow 0} \frac{\Delta y}{\Delta x} = \lim_{\Delta x \rightarrow 0} \frac{f(x_0 + \Delta x) - f(x_0)}{\Delta x}$$

常见导函数

$$C' = 0, C \text{ 为常数} \quad (x^n)' = n \cdot x^{n-1}$$

$$(a^x)' = a^x \cdot \ln a \quad (e^x)' = e^x$$

$$(\log_a x)' = \frac{1}{x \ln a} \quad (\ln x)' = \frac{1}{x}$$

$$(u \pm v)' = u' \pm v' \quad (uv)' = u'v + uv' \quad \left(\frac{u}{v}\right)' = \frac{u'v - uv'}{v^2}$$

$$y = f(g(x)) \Rightarrow y' = \frac{df}{dg} \cdot \frac{dg}{dx}$$

偏导数

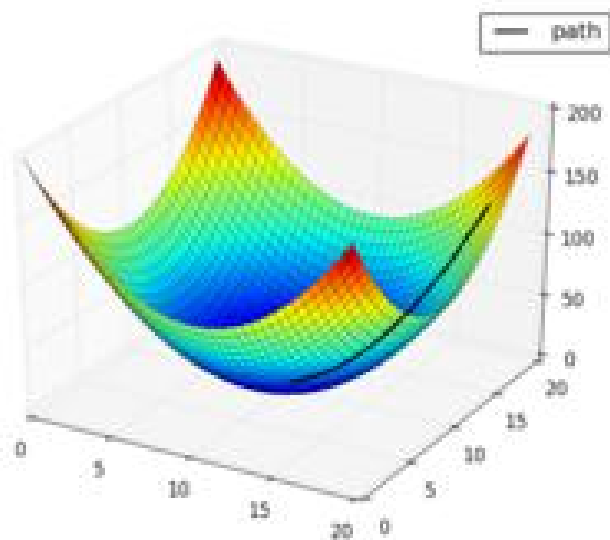
- 在一个多变量的函数中，偏导数就是关于其中一个变量的导数而保持其它变量恒定不变。假定二元函数 $z=f(x,y)$ ，点 (x_0,y_0) 是其定义域内的一个点，将 y 固定在 y_0 上，而 x 在 x_0 上增量 Δx ，相应的函数 z 有增量 $\Delta z=f(x_0+\Delta x, y_0) - f(x_0,y_0)$ ； Δz 和 Δx 的比值当 Δx 的值趋近于0的时候，如果极限存在，那么此极限值称为函数 $z=f(x,y)$ 在处对 x 的偏导数(partial derivative)，记作： $f'_x(x_0,y_0)$

$$\text{对 } x \text{ 的偏导数: } \left. \frac{\partial f}{\partial x} \right|_{\substack{x=x_0 \\ y=y_0}}$$

$$\text{对 } y \text{ 的偏导数: } \left. \frac{\partial f}{\partial y} \right|_{\substack{x=x_0 \\ y=y_0}}$$

梯度

- 梯度：梯度是一个向量，表示某一函数在该点处的**方向导数**沿着该方向取的最大值，即函数在该点处沿着该方向变化最快，变化率最大(即该梯度向量的模)；当函数为一维函数的时候，梯度其实就是导数。



$$\nabla f(x_1, x_2) = \left(\frac{\partial f(x_1, x_2)}{\partial x_1}, \frac{\partial f(x_1, x_2)}{\partial x_2} \right)$$

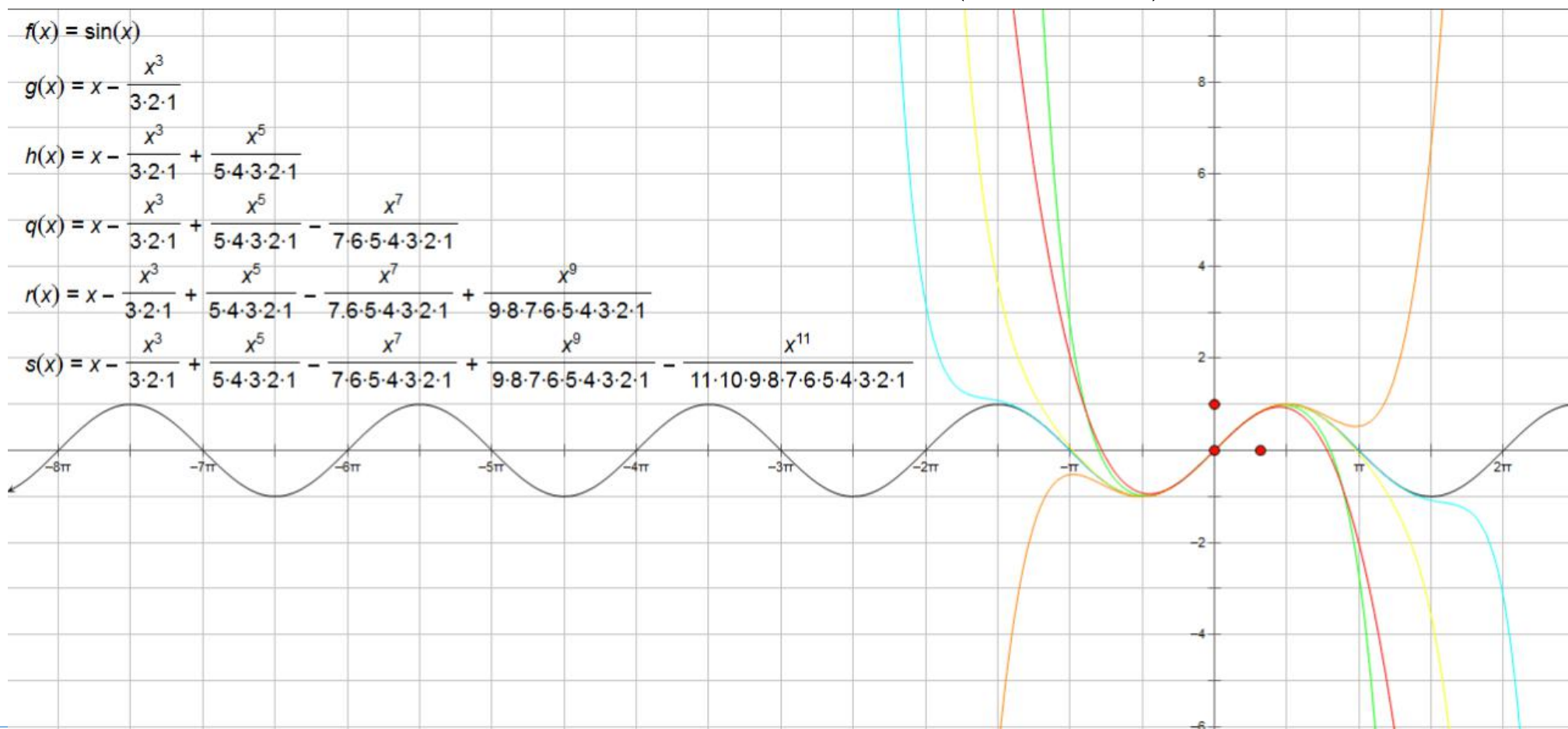
Taylor公式

- Taylor(泰勒)公式是用一个函数在某点的信息描述其附近取值的公式。如果函数足够平滑，在已知函数在某一点的各阶导数值的情况下，Taylor公式可以利用这些导数值来做系数构建一个多项式近似函数在这一点邻域中的值。
- 若函数 $f(x)$ 在包含 x_0 的某个闭区间 $[a,b]$ 上具有 n 阶函数，且在开区间 (a,b) 上具有 $n+1$ 阶函数，则对闭区间 $[a,b]$ 上任意一点 x ，有Taylor公式如下： $f^{(n)}(x)$ 表示 $f(x)$ 的 n 阶导数， $R_n(x)$ 是Taylor公式的余项，是 $(x-x_0)^n$ 的高阶无穷小

$$f(x) = \frac{f(x_0)}{0!} + \frac{f'(x_0)}{1!}(x-x_0) + \frac{f''(x_0)}{2!}(x-x_0)^2 + \dots + \frac{f^{(n)}(x_0)}{n!}(x-x_0)^n + R_n(x)$$

Taylor公式应用

$$\sin(x) = x - \frac{x^3}{3!} + \frac{x^5}{5!} + \dots + (-1)^{m-1} \frac{x^{2m-1}}{(2m-1)!} + R_{2m-1}(x)$$



Taylor公式应用2

■ 计算近似值 $e = \lim_{x \rightarrow \infty} \left(1 + \frac{1}{n}\right)^n$, 并估计误差值.

$$y = e^x \Rightarrow y' = y = e^x$$

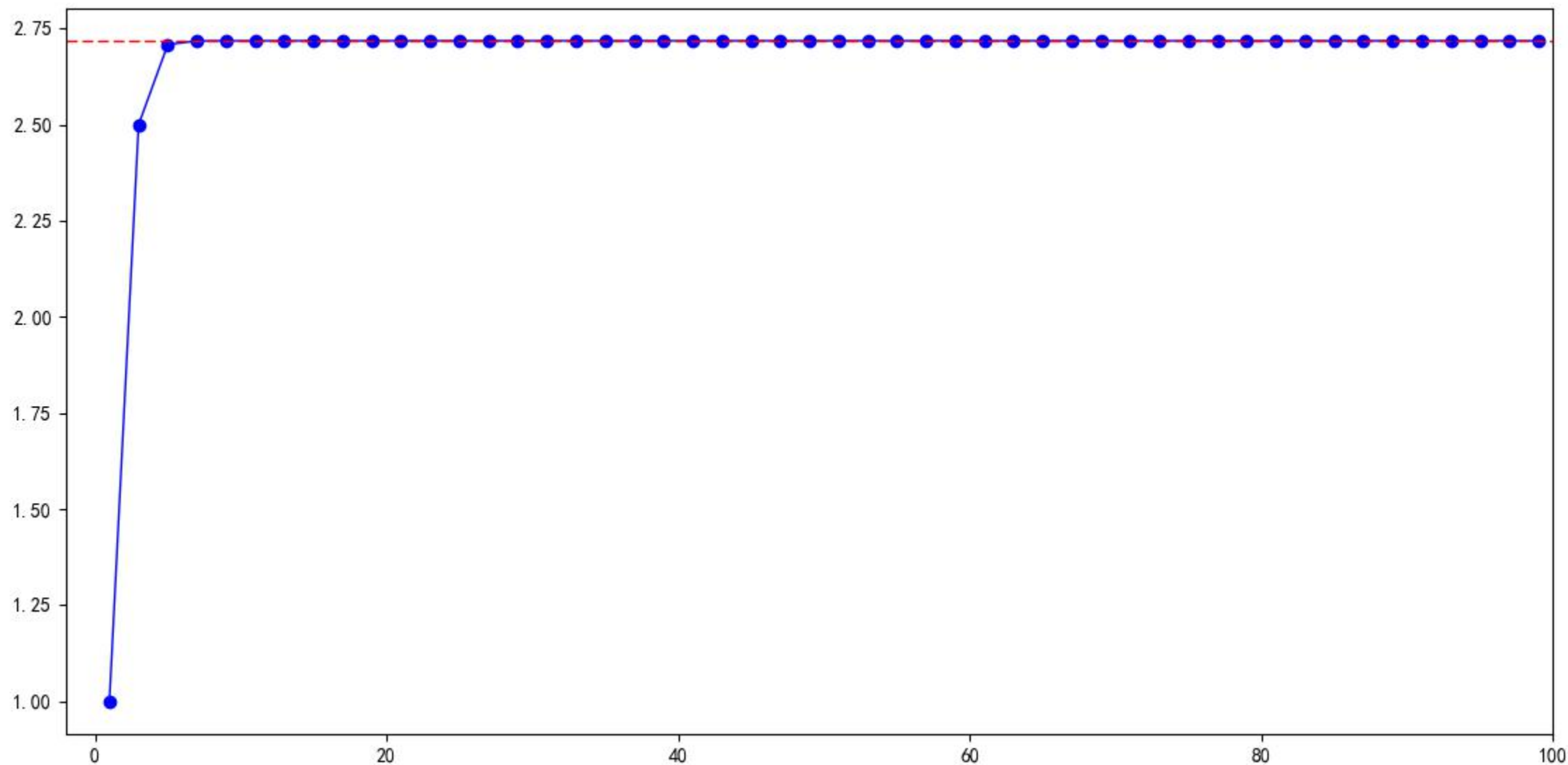
$$e^x \approx \sum_{k=0}^n \frac{e^{x_0}}{k!} (x - x_0)^k \xRightarrow{\text{令 } x_0=0} e^x \approx 1 + x + \frac{x^2}{2!} + \dots + \frac{x^n}{n!}$$

$$\xRightarrow{\text{令 } x=1} e \approx 1 + 1 + \frac{1}{2!} + \frac{1}{3!} + \dots + \frac{1}{n!}$$

$$\xRightarrow{\text{令 } n=10} e \approx 2.7182815$$

$$\delta = |R_{10}| = \frac{1}{11!} + \frac{1}{12!} + \dots = \frac{1}{11!} \left(1 + \frac{1}{12} + \frac{1}{12 \cdot 13} + \dots\right) < \frac{1}{11!} \left(1 + \frac{1}{12} + \frac{1}{12^2} + \dots\right) = \frac{12}{11 \cdot 11!} = 2.73 \cdot 10^{-8}$$

Taylor公式应用



古典概率

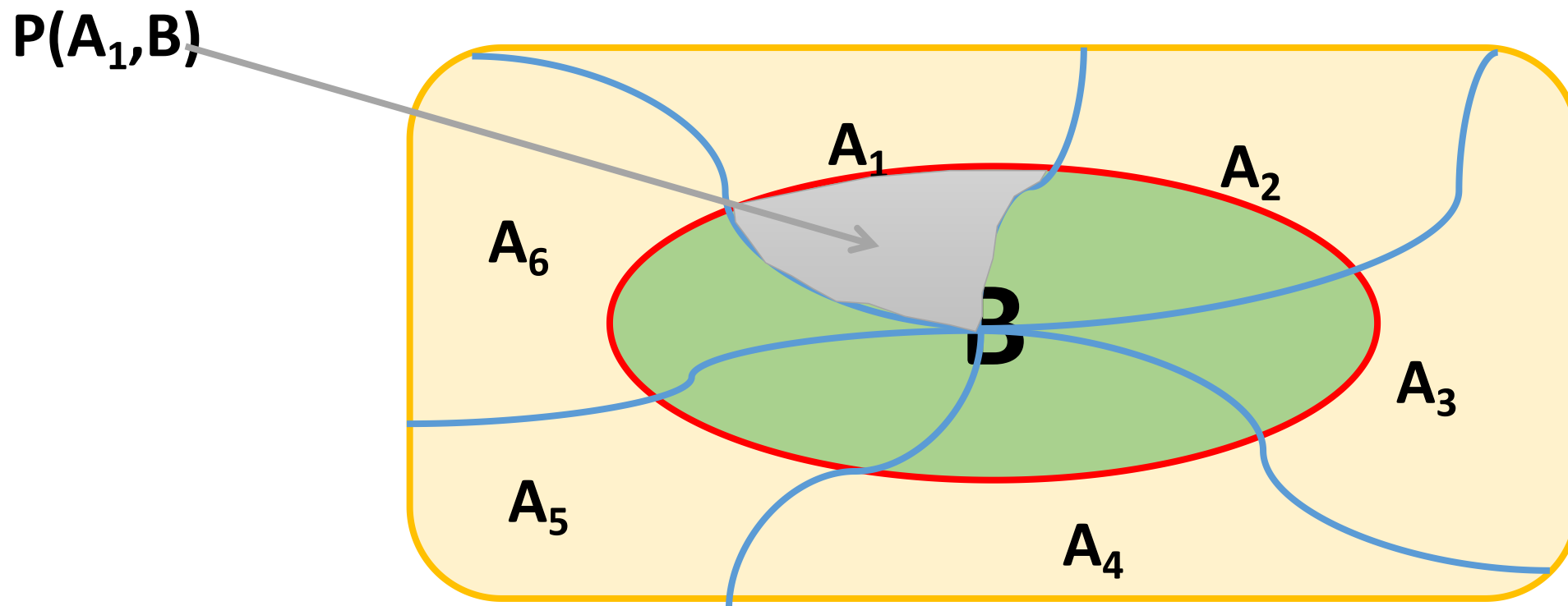
- 概率是以假设为基础的，即假定随机现象所发生的事件是有限的、互不相容的，而且每个基本事件发生的可能性相等。一般来讲，如果在全部可能出现的基本事件范围内构成事件A的基本事件有a个，不构成事件A的有b个，那么事件A出现的概率为：

$$P(A) = \frac{a}{a+b}$$

- 概率体现的是随机事件A发生可能的大小度量(数值)

联合概率

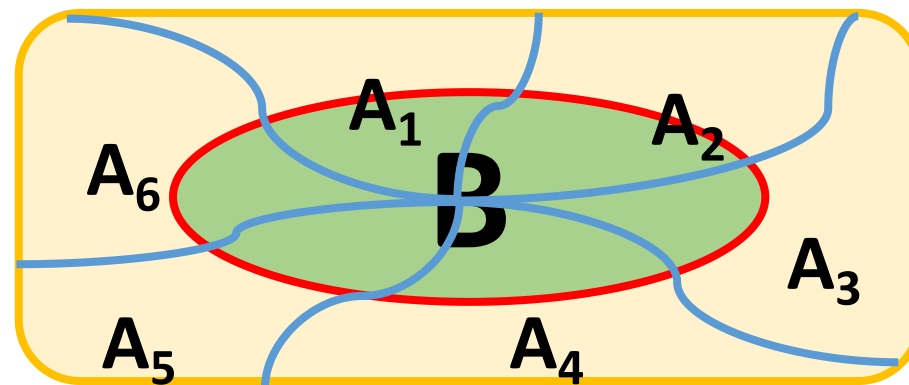
- 表示两个事件共同发生的概率，事件A和事件B的共同概率记作： $P(AB)$ 、 $P(A,B)$ 或者 $P(A \cap B)$ ，读作“事件A和事件B同时发生的概率”



条件概率

- 事件A在另外一个事件B已经发生的条件下的发生概率叫做条件概率，表示为 $P(A|B)$ ，读作“在B条件下A发生的概率”，一般情况下 $P(A|B) \neq P(A)$ ，而且条件概率具有三个特性：

- ◆ 非负性
- ◆ 可列性
- ◆ 可加性



$$P(A|B) = \frac{P(A, B)}{P(B)}$$

条件概率

- 将条件概率公式由两个事件推广到任意有穷多个事件时，可以得到如下公式，假设 A_1, A_2, \dots, A_n 为 n 个任意事件($n \geq 2$)，而且 $P(A_1 A_2 \dots A_n) > 0$ ，则：

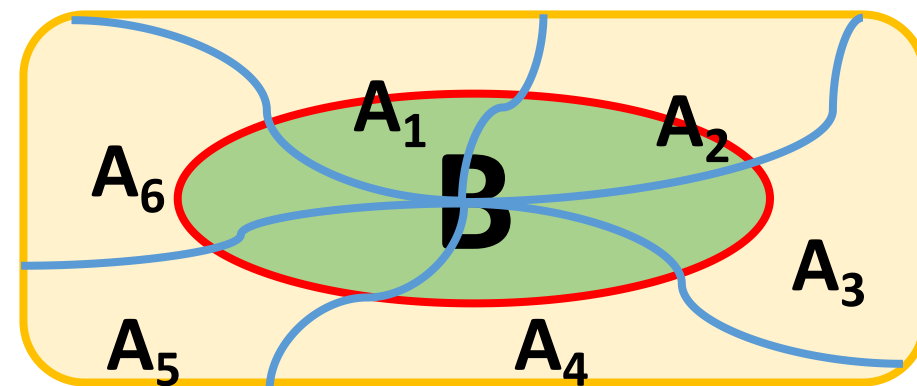
$$P(A_1 A_2 \dots A_n) = P(A_1) P(A_2 | A_1) \dots P(A_n | A_1 A_2 \dots A_{n-1})$$

全概率公式

- 样本空间 Ω 有一组事件 A_1, A_2, \dots, A_n , 如果事件组满足下列两个条件, 那么事件组称为样本空间的一个划分:

$$\forall i \neq j \in \{1, 2, \dots, n\}, A_i A_j = \phi$$

$$A_1 \cup A_2 \dots \cup A_n = \Omega$$



- 设事件 $\{A_i\}$ 是样本空间 Ω 的一个划分, 且 $P(A_i) > 0$, 那么对于任意事件 B , 全概率公式为:

$$P(B) = \sum_{i=1}^n P(A_i)P(B | A_i)$$

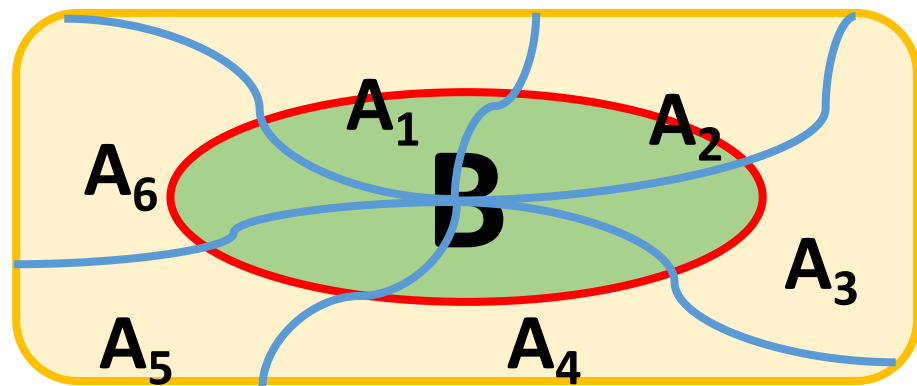
贝叶斯公式

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

IT在线教育领导品牌
EDUCATION TO CREATE A BRIGHT FUTURE

- 设 A_1, A_2, \dots, A_n 是样本空间 Ω 的一个划分，如果对任意事件 B 而言，有 $P(B) > 0$ ，那么：

$$P(A_i|B) = \frac{P(B, A_i)}{P(B)} = \frac{P(A_i) \cdot P(B|A_i)}{\sum_{j=1}^n P(A_j) \cdot P(B|A_j)}$$

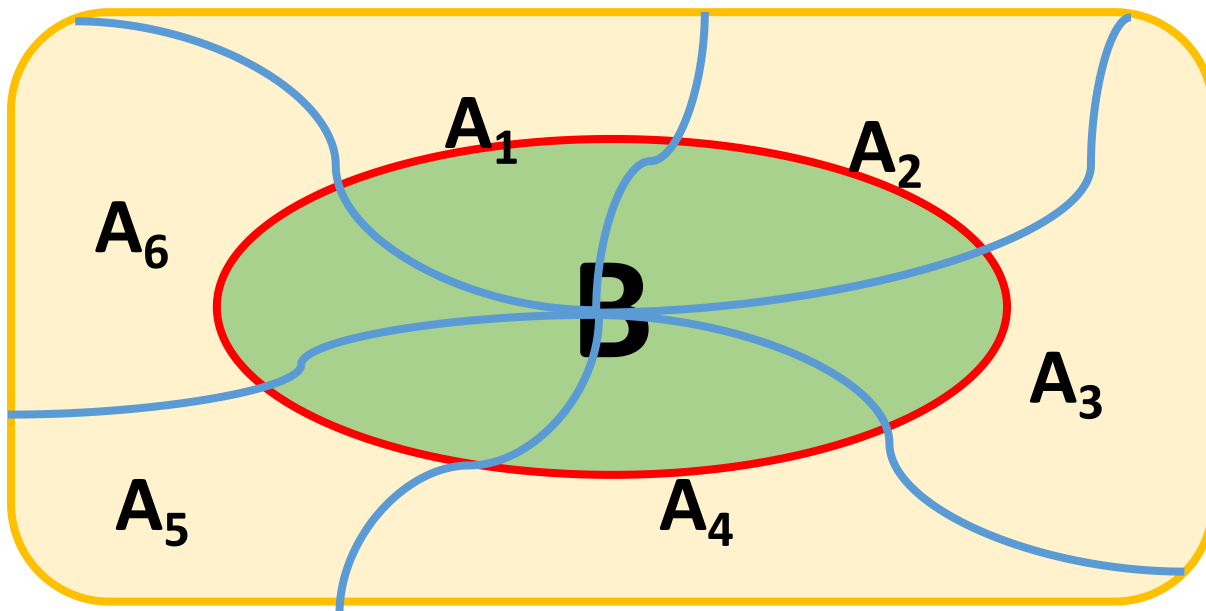


概率公式

$$P(A | B) = \frac{P(A, B)}{P(B)}$$

$$P(B) = \sum_{i=1}^n P(A_i)P(B | A_i)$$

$$P(A_i | B) = \frac{P(B | A_i)P(A_i)}{\sum_{j=1}^n P(A_j)P(B | A_j)}$$



期望

- 期望(mean)：也就是均值，是概率加权下的“平均值”，是每次可能结果的概率乘以其结果的总和，反映的实随机变量平均取值大小。常用符号 μ 表示：

- 连续性数据：

$$E(X) = \int_{-\infty}^{\infty} xf(x)dx$$

- 离散型数据：

$$E(X) = \sum_i x_i p_i$$

期望

X	2	4	6	8	10
P(x)	0.2	0.2	0.2	0.2	0.2

$$\begin{aligned}
 E(X) &= \sum_i x_i p_i \\
 &= 2 * 0.2 + 4 * 0.2 + 6 * 0.2 + 8 * 0.2 + 10 * 0.2 \\
 &= 6
 \end{aligned}$$

期望

假设C为一个常数，X和Y为两个随机变量，那么期望有以下性质：

$$E(C) = C \quad E(CX) = CE(X)$$

$$E(X + Y) = E(X) + E(Y)$$

如果X和Y相互独立，那么 $E(XY) = E(X)E(Y)$

如果 $E(XY) = E(X)E(Y)$ ，那么X和Y不相关

方差

- 方差(variance)是衡量随机变量或一组数据时离散程度的度量，是用来度量随机变量和其数学期望之间的偏离程度。即方差是衡量数据原数据和期望/均值相差的度量值。

$$Var(X) = D(X) = \sigma^2 = \frac{\sum (X - \mu)^2}{N}$$

$$D(X) = \sum_{i=1}^n p_i \cdot (x_i - \mu)^2 \qquad D(X) = \int_a^b (x - \mu)^2 f(x) dx$$

$$D(X) = E\left((X - E(X))^2\right) = E(X^2) - (E(X))^2$$

方差

X	2	4	6	8	10
P(x)	0.2	0.2	0.2	0.2	0.2

$$E(X) = 6 \quad E(X^2) = 44$$

$$\begin{aligned}
 D(X) &= E(X^2) - (E(X))^2 \\
 &= 44 - 6^2 \\
 &= 8
 \end{aligned}$$

方差

- 假设C为一个常数，X和Y是两个随机变量，那么方差有以下性质：

$$D(C) = 0 \quad D(CX) = C^2 D(X) \quad D(C + X) = D(X)$$

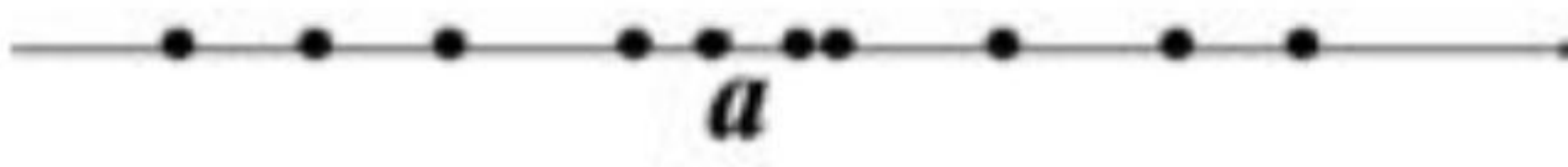
$$D(X \pm Y) = D(X) + D(Y) \pm 2Cov(X, Y)$$

$$\text{协方差 } Cov(X, Y) = E\{(X - E(X)) \cdot (Y - E(Y))\}$$

如果X和Y不相关，那么 $D(X \pm Y) = D(X) + D(Y)$

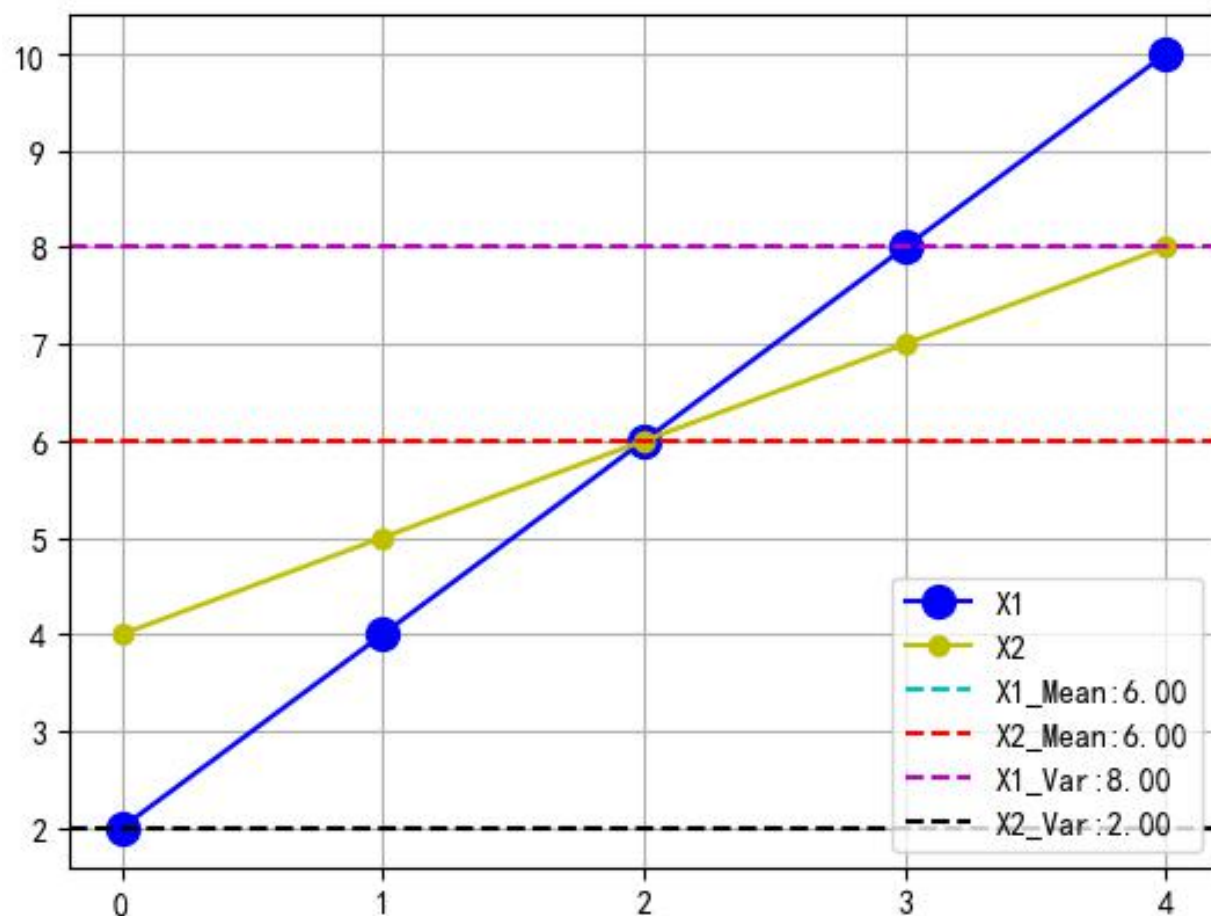
方差

- 已知某零件的真实长度为 a ，现用甲、乙两台仪器各测量10次，将测量结果 X 用坐标轴上的点表示，并且甲仪器的测量如图中的点，乙仪器的测量结果全是 a ，此时两台仪器的测量均值都是 a ，但是我们可能会认为乙机器性能更好，因为乙机器的测量值在 a 附近。由此可见，研究随机变量与其均值的偏离程度是十分必要的。



X1	2	4	6	8	10
P(x1)	0.2	0.2	0.2	0.2	0.2

X2	4	5	6	7	8
P(x2)	0.2	0.2	0.2	0.2	0.2



$$E(X_1) = 6$$

$$D(X_1) = 8$$

$$E(X_2) = 6$$

$$D(X_2) = 2$$

标准差

- 标准差(Standard Deviation)是离均值平方的算术平均数的平方根，用符号 σ 表示，其实标准差就是方差的算术平方根。
- 标准差和方差都是测量离散趋势的最重要、最常见的指标。标准差和方差的不同点在于，标准差和变量的计算单位是相同的，比方差清楚，因此在很多分析的时候使用的是标准差。

$$\sigma = \sqrt{D(X)} = \sqrt{\frac{\sum (X - \mu)^2}{N}}$$

标准差

X1	2	4	6	8	10
P(x1)	0.2	0.2	0.2	0.2	0.2

X2	4	5	6	7	8
P(x2)	0.2	0.2	0.2	0.2	0.2

$$D(X_1) = 8$$

$$\sigma_1 = \sqrt{D(X_1)} = \sqrt{8} = 2.8284$$

$$D(X_2) = 2$$

$$\sigma_2 = \sqrt{D(X_2)} = \sqrt{2} = 1.4142$$

协方差

- 协方差常用于衡量两个变量的总体误差；当两个变量相同的情况下，协方差其实就是方差。
- 如果X和Y是统计独立的，那么二者之间的协方差为零。但是如果协方差为零，那么X和Y是不相关的。

$$\begin{aligned}
 Cov(X, Y) &= E[(X - E(X)) \cdot (Y - E(Y))] \\
 &= E[XY - XE(Y) - YE(X) + E(X)E(Y)] \\
 &= E(XY) - E(X)E(Y)
 \end{aligned}$$

协方差

- 假设C为一个常数，X和Y实两个随机变量，那么协方差有性质如下所示：

$$Cov(X, Y) = Cov(Y, X)$$

$$Cov(aX, bY) = abCov(X, Y)$$

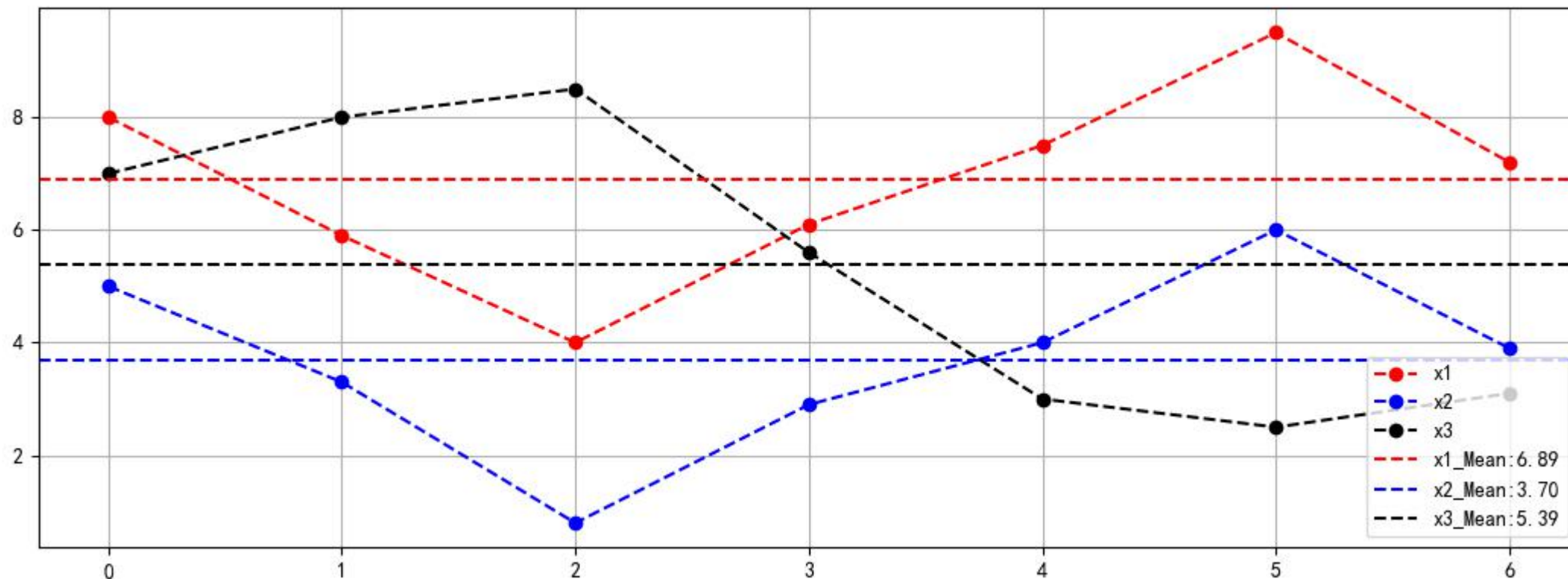
$$Cov(X_1 + X_2, Y) = Cov(X_1, Y) + Cov(X_2, Y)$$

协方差

- 协方差是两个随机变量具有相同方向变化趋势的度量：
 - ◆ 若 $\text{Cov}(X, Y) > 0$, 则X和Y的变化趋势相同；
 - ◆ 若 $\text{Cov}(X, Y) < 0$, 则X和Y的变化趋势相反；
 - ◆ 若 $\text{Cov}(X, Y) = 0$, 则X和Y不相关，也就是变化没有什么相关性

协方差

X1	8	5.9	4	6.1	7.5	9.5	7.2
X2	5	3.3	0.8	2.9	4	6	3.9
X3	7	8.0	8.5	5.6	3.0	2.5	3.1



协方差矩阵

- 对于n个随机向量($X_1, X_2, X_3, \dots, X_n$), 任意两个元素 X_i 和 X_j 都可以得到一个协方差, 从而形成一个 $n \times n$ 的矩阵, 该矩阵就叫做协方差矩阵, 协方差矩阵为对称矩阵。

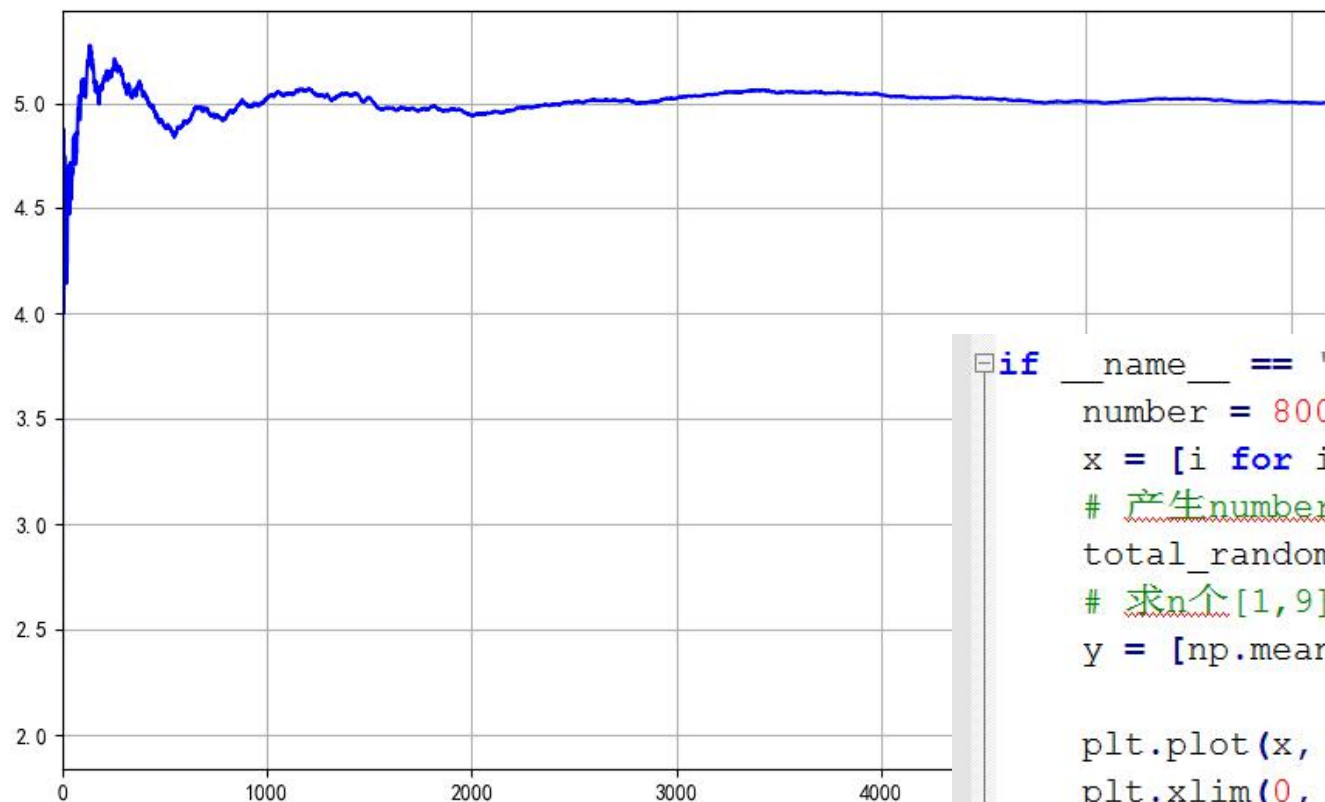
$$c_{ij} = E\{[X_i - E(X_i)][X_j - E(X_j)]\} = Cov(X_i, X_j)$$

$$C = \begin{bmatrix} c_{11} & c_{12} & \dots & c_{1n} \\ c_{21} & c_{22} & \dots & c_{2n} \\ \dots & \dots & \dots & \dots \\ c_{n1} & c_{n2} & \dots & c_{nn} \end{bmatrix}$$

大数定理

- 大数定律的意义：随着样本容量 n 的增加，样本平均数将接近于总体平均数(期望 μ)，所以在统计推断中，一般都会使用样本平均数估计总体平均数的值。
- 也就是我们会使用一部分样本的平均值来代替整体样本的期望/均值，出现偏差的可能是存在的，但是当 n 足够大的时候，偏差的可能性是非常小的，当 n 无限大的时候，这种可能性的概率基本为0。
- 大数定律的主要作用就是为使用**频率**来估计**概率**提供了理论支持；为使用部分数据来近似的模拟构建全部数据的特征提供了理论支持。

大数定理



解决中文显示问题

```
mpl.rcParams['font.sans-serif'] = [u'SimHei']
```

```
mpl.rcParams['axes.unicode_minus'] = False
```

给定随机数的种子

```
random.seed(28)
```

```
def generate_random_int(n):
```

```
    """产生n个1-9的随机数"""
```

```
    return [random.randint(1, 9) for i in range(n)]
```

```
if __name__ == '__main__':
```

```
    number = 8000
```

```
    x = [i for i in range(number + 1) if i != 0]
```

```
    # 产生number个[1,9]的随机数
```

```
    total_random_int = generate_random_int(number)
```

```
    # 求n个[1,9]的随机数的均值, n=1,2,3,4,5.....
```

```
    y = [np.mean(total_random_int[0:i + 1]) for i in range(number)]
```

```
    plt.plot(x, y, 'b-')
```

```
    plt.xlim(0, number)
```

```
    plt.grid(True)
```

```
    plt.show()
```

中心极限定理

- 中心极限定理(**Central Limit Theorem**)；假设 $\{X_n\}$ 为独立同分布的随机变量序列，并具有相同的期望 μ 和方差为 σ^2 ，则 $\{X_n\}$ 服从中心极限定理，且 Z_n 为随机序列

$\{X_n\}$ 的规范和：

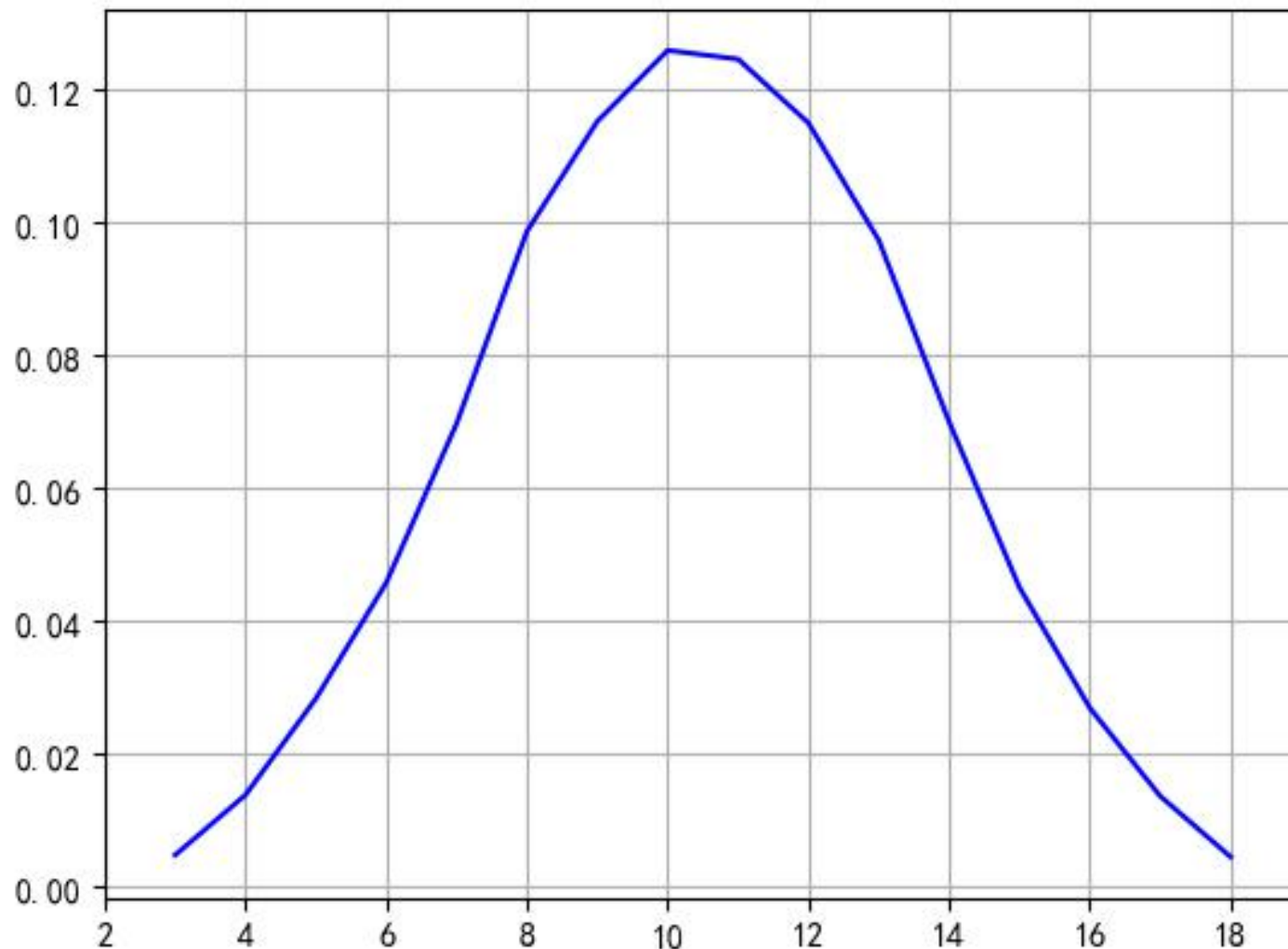
$$Y_n = X_1 + X_2 + \dots + X_n = \sum_{i=1}^n X_i \rightarrow N(n\mu, n\sigma^2)$$

$$Z_n = \frac{Y_n - E(Y_n)}{\sqrt{D(Y_n)}} = \frac{Y_n - n\mu}{\sqrt{n\sigma}} \rightarrow N(0,1)$$

- 中心极限定理就是一般在同分布的情况下，抽样样本值的规范和在总体数量趋于无穷时的极限分布近似于正态分布。

中心极限定理

- 随机的抛六面的骰子，计算三次的点数的和，三次点数的和其实就是一个事件A，现在问事件A发生的概率以及事件A所属的分布是什么？



中心极限定理

```
u"""
```

随机的抛六面的骰子，计算三次的点数的和，三次点数的和其实就是一个事件A

==> 事件A的发生属于什么分布??

==> $A = x_1 + x_2 + x_3$ ，其中 x_1 、 x_2 、 x_3 是分别三次的抛骰子的点数

根据中心极限定理，由于 x_1 、 x_2 、 x_3 属于独立同分布的，所以说最终的事件A属于高斯分布

```
"""
```

```
def generate_random_int():
```

```
    """随机产生一个[1, 6]的数字，表示的是一个六面骰子的结果"""
```

```
    return random.randint(1, 6)
```

```
def generate_mean(n):
```

```
    """计算返回n次抛六面骰子的和结果"""
```

```
    return np.sum([generate_random_int() for i in range(n)])
```

```
if __name__ == '__main__':
```

```
    # 进行A事件多少次
```

```
    number1 = 10000000
```

```
    # 表示每次A事件抛几次骰子
```

```
    number2 = 3
```

```
    # 进行number1次事件A的操作，每次事件A都进行number2次抛骰子
```

```
    keys = [generate_mean(number2) for i in range(number1)]
```

```
    # 统计每个和数字出现的次数，eg: 和为3的出现多少次、和为10出现多少次....
```

```
    result = {}
```

```
    for key in keys:
```

```
        count = 1
```

```
        if key in result:
```

```
            count += result[key]
```

```
        result[key] = count
```

```
    # 获取x和y
```

```
    x = sorted(np.unique(list(result.keys())))
```

```
    y = []
```

```
    for key in x:
```

```
        # 将出现的次数进行一个百分比的计算
```

```
        y.append(result[key] / number1)
```

```
    # 画图
```

```
    plt.plot(x, y, 'b-')
```

```
    plt.xlim(x[0] - 1, x[-1] + 1)
```

```
    plt.grid(True)
```

```
    plt.show()
```

最大似然估计

- **最大似然法(Maximum Likelihood Estimation, MLE)**也称为最大概似估计、极大似然估计，是一种具有理论性的参数估计方法。基本思想是：当从模型总体随机抽取 n 组样本观测值后，最合理的参数估计量应该使得从模型中抽取该 n 组样本观测值的概率最大；一般步骤如下：
 - ◆ 1. 写出似然函数；
 - ◆ 2. 对似然函数取对数，并整理；
 - ◆ 3. 求导数；
 - ◆ 4. 解似然方程

最大似然估计

- 设总体分布为 $f(x, \theta)$, $\{X_n\}$ 为该总体采样得到的样本。因为随机序列 $\{X_n\}$ 独立同分布，则它们的联合密度函数为：

$$L(x_1, x_2, \dots, x_n; \theta_1, \theta_2, \dots, \theta_n) = \prod_{i=1}^n f(x_i; \theta_1, \theta_2, \dots, \theta_n)$$

- 这里 θ 被看做固定但是未知的参数，反过来，因为样本已经存在，可以看做 $\{X_n\}$ 是固定的， $L(x, \theta)$ 是关于 θ 的函数，即似然函数；
- 求参数 θ 的值，使得似然函数取最大值，这种方法叫做**最大似然估计法**。

最大似然估计

- 若给定一组样本 $\{X_n\}$ ，已知随机样本符合高斯分布 $N(\mu, \sigma^2)$ ，试估计 σ 和 μ 的值

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} \cdot e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad L(x) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \cdot e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$l(x) = \log(L(x)) = \log \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \cdot e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$= \sum_{i=1}^n \log \left(\frac{1}{\sqrt{2\pi}\sigma} \cdot e^{-\frac{(x-\mu)^2}{2\sigma^2}} \right) = \sum_{i=1}^n \log \left(\frac{1}{\sqrt{2\pi}\sigma} \right) - \sum_{i=1}^n \frac{(x-\mu)^2}{2\sigma^2} = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_i (x-\mu)^2$$

最大似然估计

$$l(x) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_i (x_i - \mu)^2$$

$$= -\frac{n}{2} \log(2\pi) - \left[n \log(\sigma) + \frac{1}{2\sigma^2} \sum_i (x_i - \mu)^2 \right] = -\frac{n}{2} \log(2\pi) - l_2(x)$$

- 要求似然函数 $l(x)$ 最大，即 $l(x)$ 求极值即可，将似然函数对参数 μ 和 σ 分别求偏导数

$$\frac{dl(x)}{d\mu} = -\frac{1}{\sigma^2} \sum_i (x_i - \mu) = 0 \Rightarrow \mu = \frac{1}{n} \sum_i x_i$$

$$\frac{dl_2(x)}{d\sigma} = \frac{\pi}{\sigma} - \frac{\sum_i (x_i - \mu)^2}{\sigma^3} = 0 \Rightarrow \sigma^2 = \frac{1}{n} \sum_i (x_i - \mu)^2$$

向量的运算

- 设两向量为： $\vec{a} = (x_1, y_1)$ ， $\vec{b} = (x_2, y_2)$
- 向量的加法/减法满足平行四边形法则和三角形法则

$$\vec{a} + \vec{b} = (x_1 + x_2, y_1 + y_2)$$

$$\vec{a} - \vec{b} = (x_1 - x_2, y_1 - y_2)$$

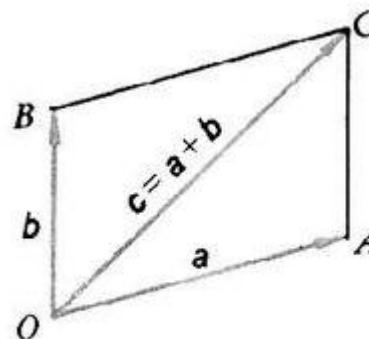


图4 向量的加法

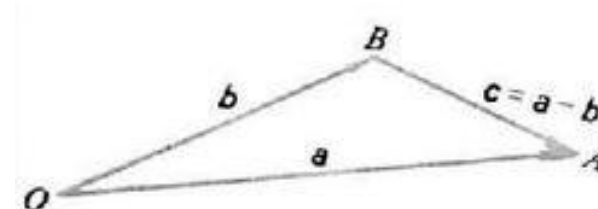


图5 向量的减法

- 数乘：实数 λ 和向量 a 的叉乘乘积还是一个向量，记作 λa ，且 $|\lambda a| = |\lambda| |a|$ ；数乘的几何意义是将向量 a 进行伸长或者压缩操作

$$\lambda \vec{a} = (\lambda x_1, \lambda y_1)$$

向量的运算

■ 设两向量为： $\vec{a} = (x_1, y_1)$ ， $\vec{b} = (x_2, y_2)$ ，并且a和b之间的夹角为： θ

■ 数量积：两个向量的数量积(内积、点积)是一个数量/实数，记作 $\vec{a} \cdot \vec{b}$

$$\vec{a} \cdot \vec{b} = |\vec{a}| * |\vec{b}| * \cos \theta$$

■ 向量积：两个向量的向量积(外积、叉积)是一个向量，记作 $\vec{a} \times \vec{b}$ ；

向量积即两个不共线非零向量所在平面的一组法向量。

$$|\vec{a} \times \vec{b}| = |\vec{a}| * |\vec{b}| * \sin \theta$$

矩阵的直观表示

- 数域F中 $m \times n$ 个数排成m行n列，并括以圆括弧(或方括弧)的数表示成为数域F上的矩阵，通常用大写字母记作A或者 $A_{m \times n}$ ，有时也记作 $A = (a_{ij})_{m \times n} (i=1, 2, \dots, m; j=1, 2, \dots, n)$ ，其中 a_{ij} 表示矩阵A的第i行的第j列元素，当F为实数域R时，A叫做实矩阵，当F为复数域C时，A叫做复矩阵。

$$A = \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \cdots & \cdots & \cdots & \cdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{pmatrix}$$

矩阵的加减法

- 矩阵的加法与减法要求进行操作的两个矩阵A和B具有相同的阶，假设A为m*n阶矩阵，B为m*n阶矩阵，那么C=A ± B也是m*n阶的矩阵，并且矩阵C的元素满足： $c_{ij} = a_{ij} \pm b_{ij}$

$$A = \begin{Bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \dots & \dots & \dots & \dots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{Bmatrix} \quad B = \begin{Bmatrix} b_{11} & b_{12} & \dots & b_{1n} \\ b_{21} & b_{22} & \dots & b_{2n} \\ \dots & \dots & \dots & \dots \\ b_{m1} & b_{m2} & \dots & b_{mn} \end{Bmatrix} \quad C = A \pm B = \begin{Bmatrix} a_{11} \pm b_{11} & a_{12} \pm b_{12} & \dots & a_{1n} \pm b_{1n} \\ a_{21} \pm b_{21} & a_{22} \pm b_{22} & \dots & a_{2n} \pm b_{2n} \\ \dots & \dots & \dots & \dots \\ a_{m1} \pm b_{m1} & a_{m2} \pm b_{m2} & \dots & a_{mn} \pm b_{mn} \end{Bmatrix}$$

■ 记 $A_{m \times n} = (a_{ij})$ $B_{m \times n} = (b_{ij})$

加法： $A_{m \times n} + B_{m \times n} = (a_{ij} + b_{ij})$

$$\begin{bmatrix} 1 & 2 \\ 5 & 6 \end{bmatrix} + \begin{bmatrix} 3 & 4 \\ 7 & 8 \end{bmatrix} = \begin{bmatrix} 4 & 6 \\ 12 & 14 \end{bmatrix}$$

减法：

$$A_{m \times n} - B_{m \times n} = (a_{ij} - b_{ij})$$

$$\begin{bmatrix} 6 & 8 \\ 11 & 14 \end{bmatrix} - \begin{bmatrix} 3 & 4 \\ 7 & 8 \end{bmatrix} = \begin{bmatrix} 3 & 4 \\ 4 & 6 \end{bmatrix}$$

运算律：

交换律 $A + B = B + A$

结合律 $(A + B) + C = A + (B + C)$

矩阵与数的乘法

- 数乘：将数 λ 与矩阵 A 相乘，就是将数 λ 与矩阵 A 中的每一个元素相乘，记作 λA ；结果 $C=\lambda A$ ，并且 C 中的元素满足： $c_{ij} = \lambda a_{ij}$

$$A = \begin{Bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \cdots & \cdots & \cdots & \cdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{Bmatrix}$$

$$\lambda A = \begin{Bmatrix} \lambda a_{11} & \lambda a_{12} & \cdots & \lambda a_{1n} \\ \lambda a_{21} & \lambda a_{22} & \cdots & \lambda a_{2n} \\ \cdots & \cdots & \cdots & \cdots \\ \lambda a_{m1} & \lambda a_{m2} & \cdots & \lambda a_{mn} \end{Bmatrix}$$

矩阵与数的乘法

■ 数乘： $kA = (ka_{ij})$

$$3 \times \begin{bmatrix} 1 & 2 \\ 5 & 6 \end{bmatrix} = \begin{bmatrix} 3 \times 1 & 3 \times 2 \\ 3 \times 5 & 3 \times 6 \end{bmatrix} = \begin{bmatrix} 3 & 6 \\ 15 & 18 \end{bmatrix}$$

■ 数乘运算律

结合律 $(\lambda\mu)A = \lambda(\mu A)$

分配律 $(\lambda + \mu)A = \lambda A + \mu A$

$$\lambda(A + B) = \lambda A + \lambda B$$

矩阵与向量的乘法

- 假设A为m*n阶矩阵，x为n*1的列向量，则Ax为m*1的列向量，记作： $\vec{y} = A \vec{x}$

$$A = \begin{Bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \dots & \dots & \dots & \dots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{Bmatrix} \quad \vec{x} = \begin{pmatrix} x_1 \\ x_2 \\ \dots \\ x_n \end{pmatrix} \quad \vec{y} = A \vec{x} = \begin{pmatrix} y_1 \\ y_2 \\ \dots \\ y_m \end{pmatrix}$$

$$y_i = \sum_{j=1}^n a_{ij} x_j$$

矩阵与矩阵的乘法

- 矩阵的乘法仅当第一个矩阵A的列数和第二个矩阵B的行数相等时才能够定义，假设A为m*s阶矩阵，B为s*n阶矩阵，那么C=A*B是m*n阶矩阵，并且矩阵C中的元素满足：
$$c_{ij} = \sum_{k=1}^s a_{ik} b_{kj}$$

$$A = \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1s} \\ a_{21} & a_{22} & \cdots & a_{2s} \\ \cdots & \cdots & \cdots & \cdots \\ a_{m1} & a_{m2} & \cdots & a_{ms} \end{pmatrix} \quad B = \begin{pmatrix} b_{11} & b_{12} & \cdots & b_{1n} \\ b_{21} & b_{22} & \cdots & b_{2n} \\ \cdots & \cdots & \cdots & \cdots \\ b_{s1} & b_{s2} & \cdots & b_{sn} \end{pmatrix} \quad C = A * B = \begin{pmatrix} c_{11} & c_{12} & \cdots & c_{1n} \\ c_{21} & c_{22} & \cdots & c_{2n} \\ \cdots & \cdots & \cdots & \cdots \\ c_{m1} & c_{m2} & \cdots & c_{mn} \end{pmatrix}$$

矩阵与矩阵的乘法

■ 例：

$$A = \begin{bmatrix} 1 & 0 & -1 & 2 \\ -1 & 1 & 3 & 0 \\ 0 & 5 & -1 & 4 \end{bmatrix} \quad B = \begin{bmatrix} 0 & 3 & 4 \\ 1 & 2 & 1 \\ 3 & 1 & -1 \\ -1 & 2 & 1 \end{bmatrix} \quad \text{求 } C=AB$$

$$C = AB = \begin{bmatrix} 1 & 0 & -1 & 2 \\ -1 & 1 & 3 & 0 \\ 0 & 5 & -1 & 4 \end{bmatrix} \begin{bmatrix} 0 & 3 & 4 \\ 1 & 2 & 1 \\ 3 & 1 & -1 \\ -1 & 2 & 1 \end{bmatrix} = \begin{bmatrix} -5 & 6 & 7 \\ 10 & 2 & -6 \\ -2 & 17 & 10 \end{bmatrix}$$

$$(AB)C = A(BC) \quad (A+B)C = AC + BC \quad C(A+B) = CA + CB$$

矩阵的转置

- 矩阵的转置：把矩阵A的行和列互相交换所产生的矩阵称为A的转置矩阵，这一过程叫做矩阵的转置。 使用 A^T 表示A的转置

$$A = \begin{Bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \cdots & \cdots & \cdots & \cdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{Bmatrix} \quad A^T = \begin{Bmatrix} a_{11} & a_{21} & \cdots & a_{m1} \\ a_{12} & a_{22} & \cdots & a_{m2} \\ \cdots & \cdots & \cdots & \cdots \\ a_{1n} & a_{2n} & \cdots & a_{mn} \end{Bmatrix}$$

矩阵的转置

■ 例： $A = \begin{bmatrix} 1 & 2 \\ 5 & 6 \end{bmatrix}$, A的转置为 $A^T = \begin{bmatrix} 1 & 5 \\ 2 & 6 \end{bmatrix}$

■ 转置的运算性质：

$$(A^T)^T = A$$

$$(\lambda A)^T = \lambda A^T$$

$$(AB)^T = B^T A^T$$

$$(A + B)^T = A^T + B^T$$

QR分解

- QR分解是将矩阵分解为一个正交矩阵与上三角矩阵的乘积

$$\begin{array}{ccc}
 \begin{array}{|c|} \hline A \\ \hline \end{array} & = & \begin{array}{|c|} \hline Q \\ \hline \end{array} \begin{array}{|c|} \hline \begin{array}{c} R \\ 0 \end{array} \\ \hline \end{array} \\
 m \times n & & m \times m \quad m \times n
 \end{array}$$

- 这其中， Q 为正交矩阵， $Q^T Q = I$ ， R 为上三角矩阵。
- 实际中，QR分解经常被用来解线性最小二乘问题。

SVD分解

- 奇异值分解(Singular Value Decomposition)是一种重要的矩阵分解方法，可以看做是对称方阵在任意矩阵上的推广。
- 假设A为一个 $m \times n$ 阶实矩阵，则存在一个分解使得：

$$A_{m \times n} = U_{m \times m} \Sigma_{m \times n} V_{n \times n}^T$$

- ◆ 通常将奇异值由大到小排列，这样 Σ 便能由A唯一确定了。

向量的导数

■ A为m*n的矩阵，x为n*1的列向量，则Ax为m*1的列向量，记作 $\vec{y} = A \cdot \vec{x}$

$$A = \begin{Bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \dots & \dots & \dots & \dots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{Bmatrix} \quad \vec{x} = \begin{pmatrix} x_1 \\ x_2 \\ \dots \\ x_n \end{pmatrix} \quad \vec{y} = \begin{pmatrix} a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n \\ a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n \\ \dots \\ a_{m1}x_1 + a_{m2}x_2 + \dots + a_{mn}x_n \end{pmatrix}$$

$$\frac{\partial \vec{y}}{\partial \vec{x}} = \frac{\partial A \vec{x}}{\partial \vec{x}} = \begin{pmatrix} a_{11} & a_{21} & \dots & a_{m1} \\ a_{12} & a_{22} & \dots & a_{m2} \\ \dots & \dots & \dots & \dots \\ a_{1n} & a_{2n} & \dots & a_{mn} \end{pmatrix} = A^T$$

向量的导数

■ 向量偏导公式

$$\frac{\partial A\vec{x}}{\partial \vec{x}} = A^T$$

$$\frac{\partial A\vec{x}}{\partial \vec{x}^T} = A$$

$$\frac{\partial (\vec{x}^T A)}{\partial \vec{x}} = A$$

标量对向量的导数

■ A为n*n的矩阵，x为n*1的列向量，记 $y = \vec{x}^T \cdot A \cdot \vec{x}$

■ 同理可得：
$$\frac{\partial y}{\partial \vec{x}} = \frac{\partial (\vec{x}^T \cdot A \cdot \vec{x})}{\partial \vec{x}} = (A^T + A) \cdot \vec{x}$$

■ 若A为对称矩阵，则有
$$\frac{\partial (\vec{x}^T A \vec{x})}{\partial \vec{x}} = 2A\vec{x}$$

标量对向量的导数

$$A = \begin{Bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \dots & \dots & \dots & \dots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{Bmatrix} \quad \vec{x} = \begin{pmatrix} x_1 \\ x_2 \\ \dots \\ x_n \end{pmatrix}$$

$$\vec{x}^T \cdot A \cdot \vec{x} = (x_1, x_2, \dots, x_n) \left(\sum_{j=1}^n a_{1j} x_j, \sum_{j=1}^n a_{2j} x_j, \dots, \sum_{j=1}^n a_{nj} x_j \right)^T = \sum_{i=1}^n \sum_{j=1}^n a_{ij} x_i x_j$$

$$\frac{\partial \left(\vec{x}^T \cdot A \cdot \vec{x} \right)}{\partial \vec{x}} = \left(\sum_{j=1}^n a_{ij} x_j \right) + \left(\sum_{j=1}^n a_{ji} x_j \right) = \sum_{j=1}^n (a_{ij} + a_{ji}) x_j$$

标量对方阵的导数

■ A为n*n的矩阵，|A|为A的行列式，计算 $\frac{\partial |A|}{\partial A}$

$$\forall 1 \leq i \leq n, |A| = \sum_{j=1}^n a_{ij} \cdot (-1)^{i+j} M_{ij}$$

$$\frac{\partial |A|}{\partial a_{ij}} = \frac{\partial \left(\sum_{j=1}^n a_{ij} \cdot (-1)^{i+j} M_{ij} \right)}{\partial a_{ij}} = (-1)^{i+j} M_{ij} = A^*_{ji}$$

$$\frac{\partial |A|}{\partial A} = (A^*)^T = |A| \cdot (A^{-1})^T$$

数学知识回顾

- 常见函数
- 导数、**梯度**
- **Taylor公式**
- **联合概率、条件概率、全概率公式、贝叶斯公式**
- 期望、**方差、协方差**
- 大数定理、**中心极限定理**
- **最大似然估计(MLE)**
- **向量、矩阵的运算**
- 向量、矩阵的求导
- **SVD分解、QR分解**

梯度下降案例 $y = f(x) = x^2$

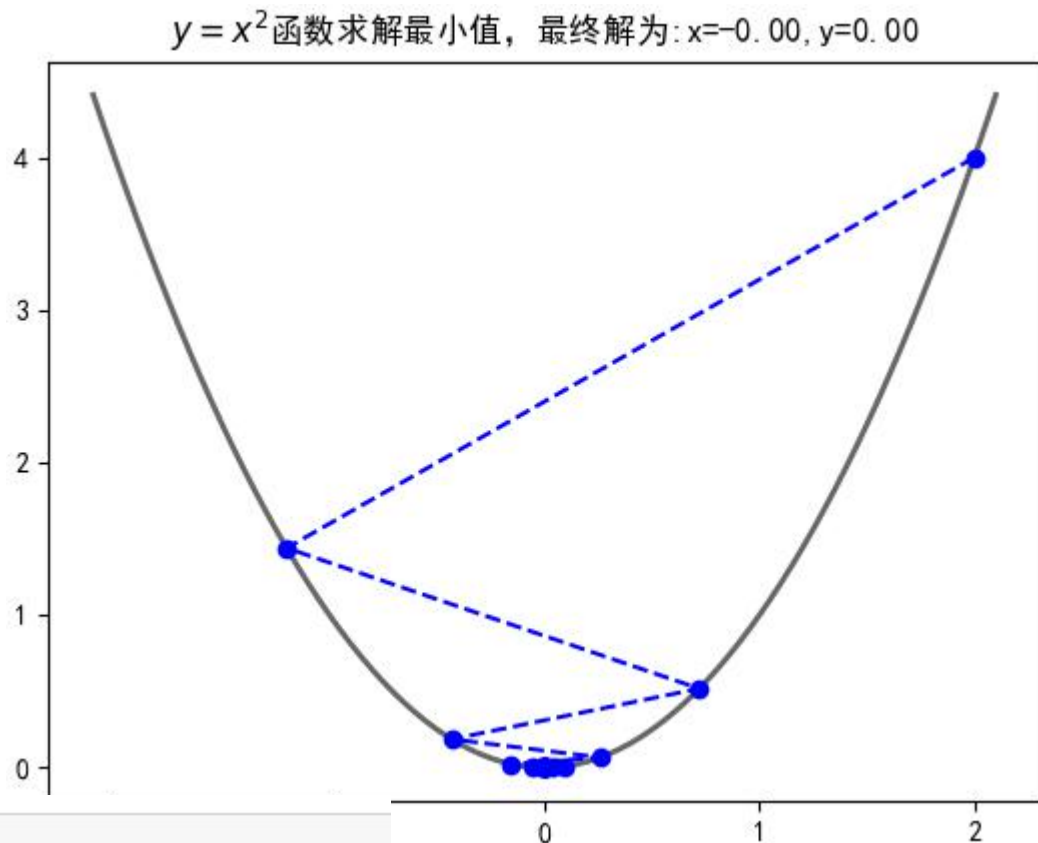
```
## 原函数
def f(x):
    return x ** 2
## 导数
def h(x):
    return 2 * x

X = []
Y = []

x = 2
step = 0.8
f_change = f(x)
f_current = f(x)
X.append(x)
Y.append(f_current)
while f_change > 1e-10:
    x = x - step * h(x)
    tmp = f(x)
    f_change = np.abs(f_current - tmp)
    f_current = tmp
    X.append(x)
    Y.append(f_current)
```

```
fig = plt.figure()
X2 = np.arange(-2.1, 2.15, 0.05)
Y2 = X2 ** 2

plt.plot(X2, Y2, '-', color='#666666', linewidth=2)
plt.plot(X, Y, 'bo—')
plt.title(u'$y=x^2$函数求解最小值, 最终解为: x=%.2f, y=%.2f'
          % (x, f_current))
plt.show()
```



梯度下降案例

$$z = f(x, y) = x^2 + y^2$$

梯度下降法求解，最终解为：x=0.00, y=0.00, z=0.00

```
## 原函数
def f(x, y):
    return x ** 2 + y ** 2
## 偏函数
def h(t):
    return 2 * t

X = []
Y = []
Z = []

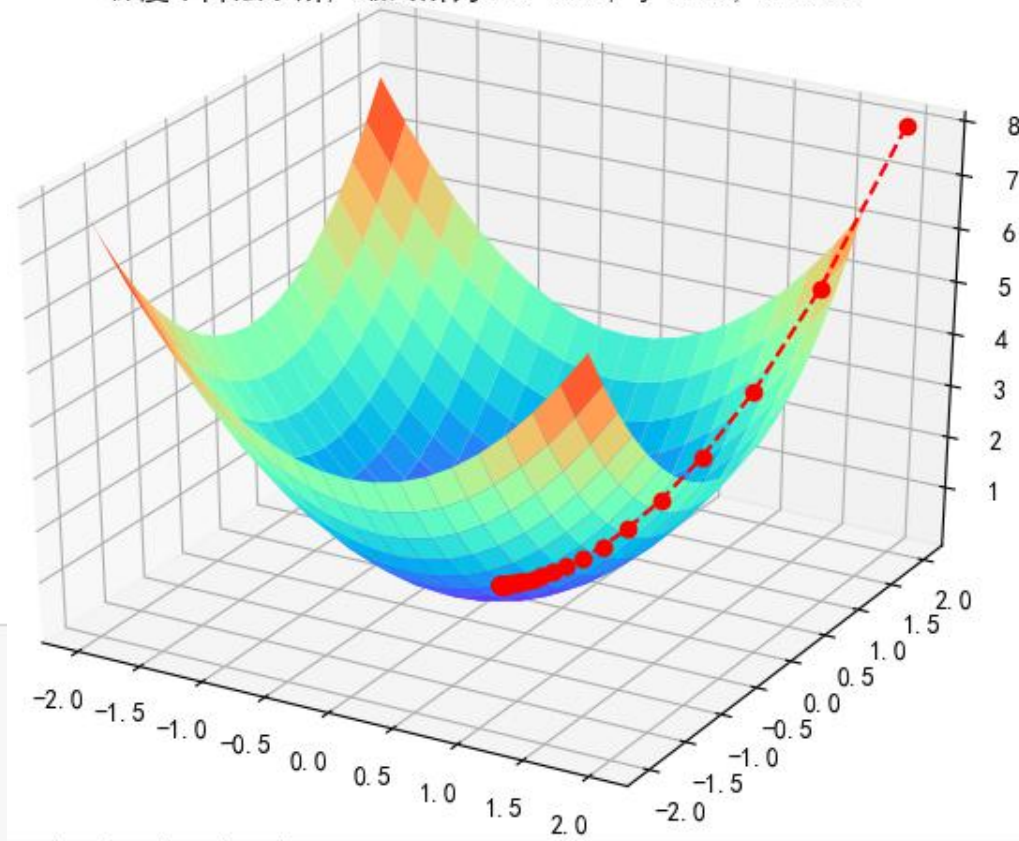
x = 2
y = 2
f_change = x ** 2 + y ** 2
f_current = f(x, y)
step = 0.1
X.append(x)
Y.append(y)
Z.append(f_current)
while f_change > 1e-10:
    x = x - step * h(x)
    y = y - step * h(y)
    f_change = f_current - f(x, y)
    f_current = f(x, y)
    X.append(x)
    Y.append(y)
    Z.append(f_current)
```

```
fig = plt.figure()
ax = Axes3D(fig)
X2 = np.arange(-2, 2, 0.2)
Y2 = np.arange(-2, 2, 0.2)
X2, Y2 = np.meshgrid(X2, Y2)
Z2 = X2 ** 2 + Y2 ** 2

ax.plot_surface(X2, Y2, Z2, rstride=1, cstride=1, cmap='rainbow')
ax.plot(X, Y, Z, 'ro—')

ax.set_title(u'梯度下降法求解，最终解为：x=%.2f, y=%.2f, z=%.2f' % (x, y, f_current))

plt.show()
```



梯度下降法

- 梯度下降法(Gradient Descent, GD)常用于求解**无约束**情况下**凸函数(Convex Function)**的**极小值**，是一种**迭代类型**的算法，因为凸函数只有一个极值点，故求解出来的极小值点就是函数的**最小值点**。

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m \left(h_{\theta}(x^{(i)}) - y^{(i)} \right)^2$$

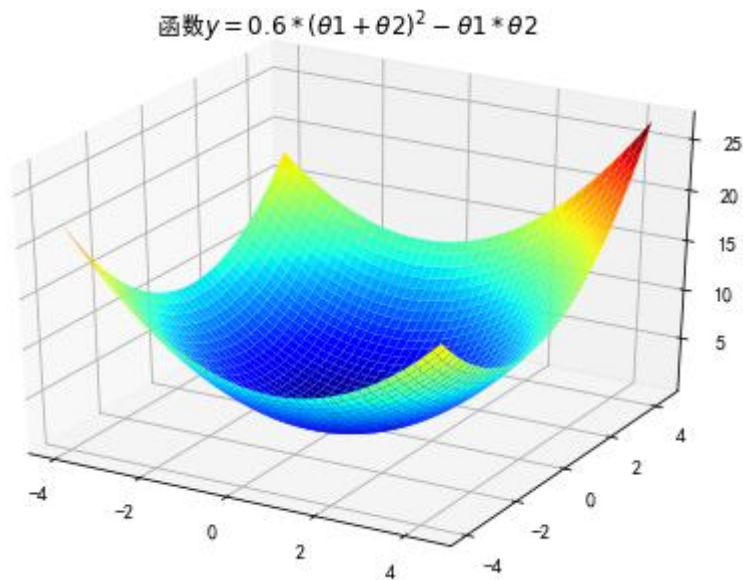
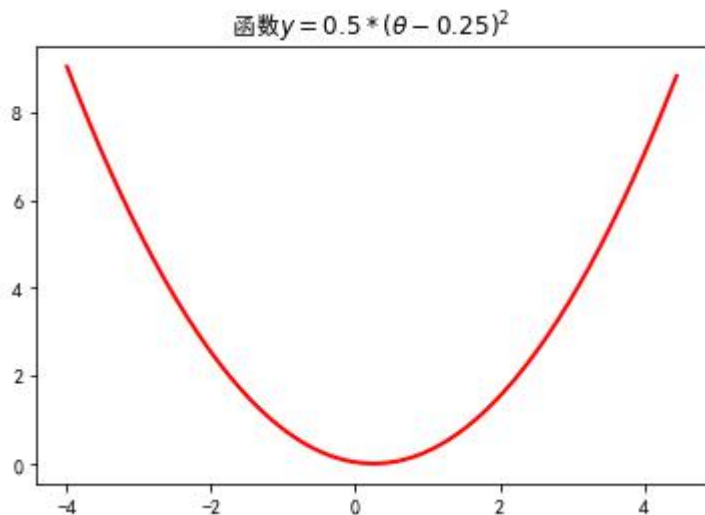
$$\theta^* = \arg \min_{\theta} J(\theta)$$

梯度下降法

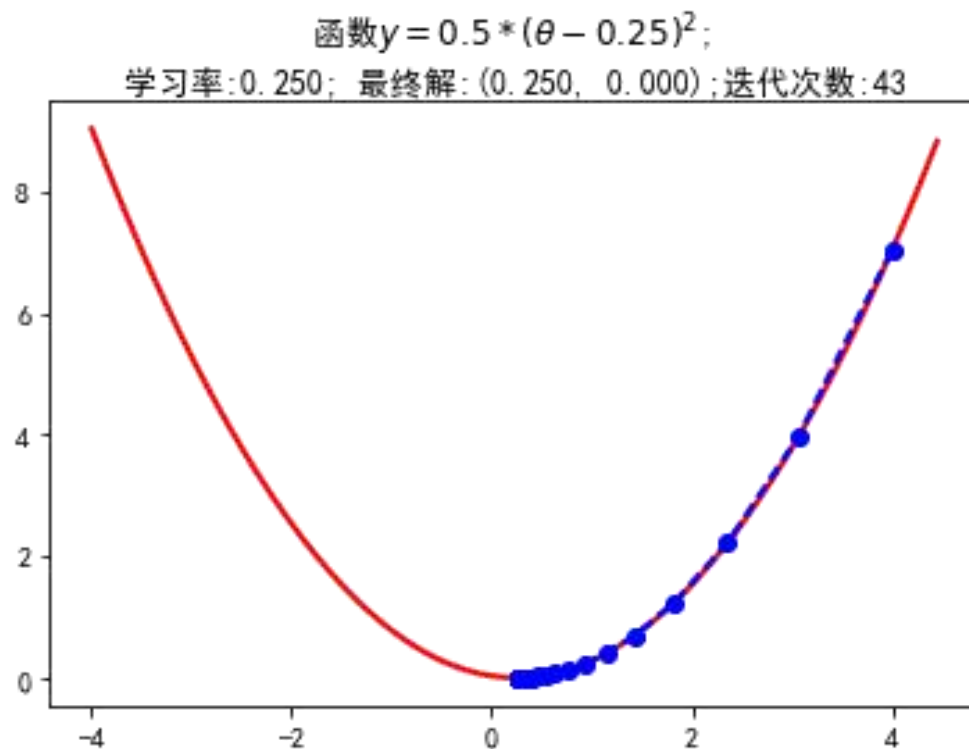
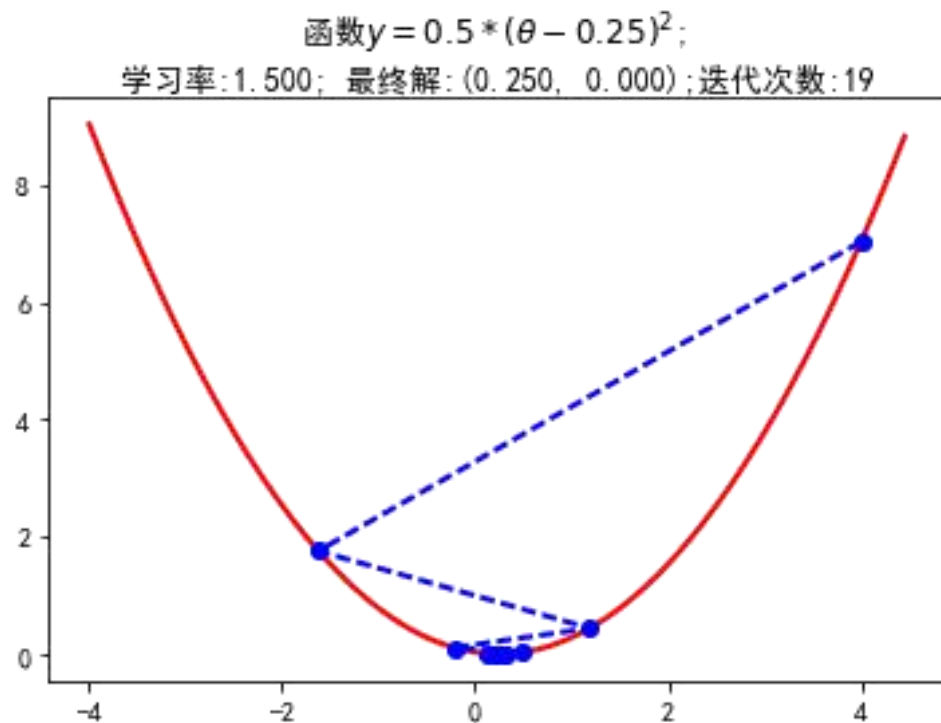
- 梯度下降法的优化思想是用当前位置负梯度方向作为搜索方向，因为该方向为当前位置的最快下降方向，所以梯度下降法也被称为“最速下降法”。梯度下降法中越接近目标值，变量变化越小。计算公式如下：

$$\theta^{k+1} = \theta^k - \alpha \nabla f(\theta^k)$$

- α 被称为**步长**或者**学习率(learning rate)**，表示自变量x每次迭代变化的大小。
- 收敛条件：当目标函数的函数值变化非常小的时候或者达到最大迭代次数的时候，就结束循环。

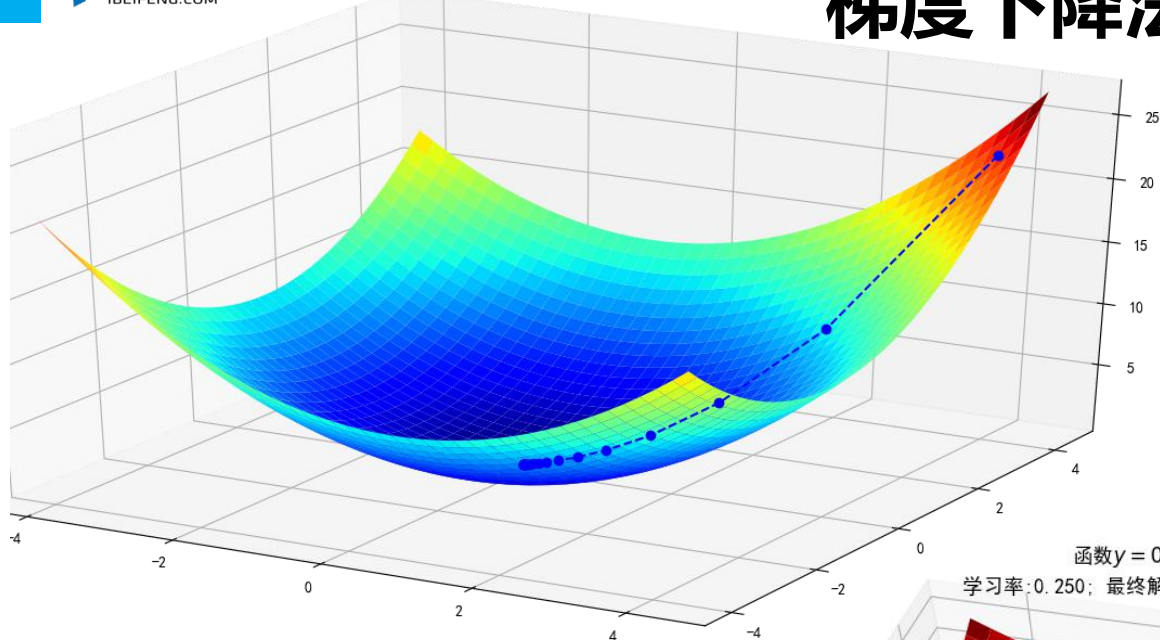


梯度下降法案例代码

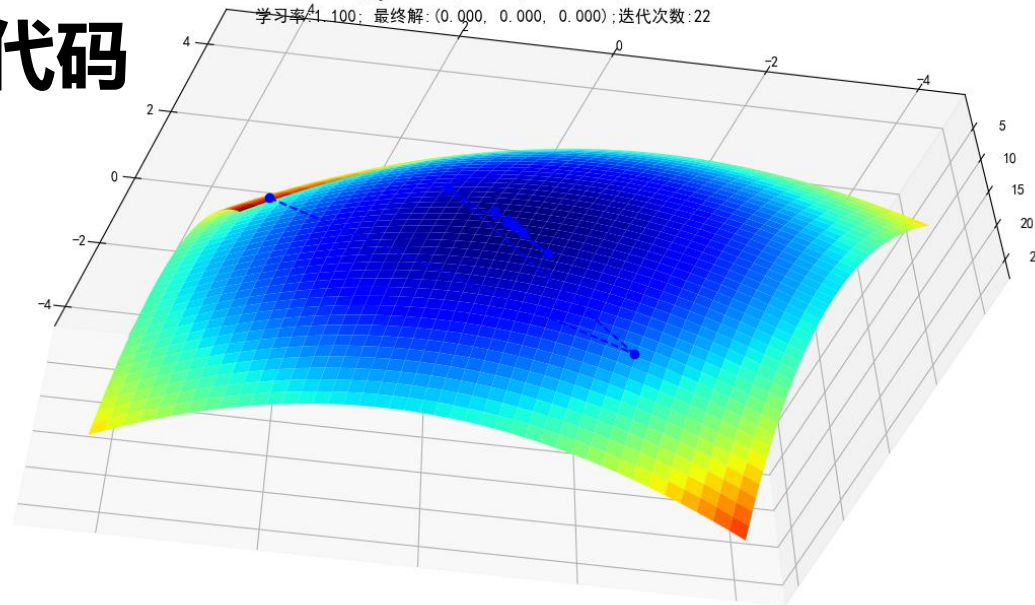


梯度下降法案例代码

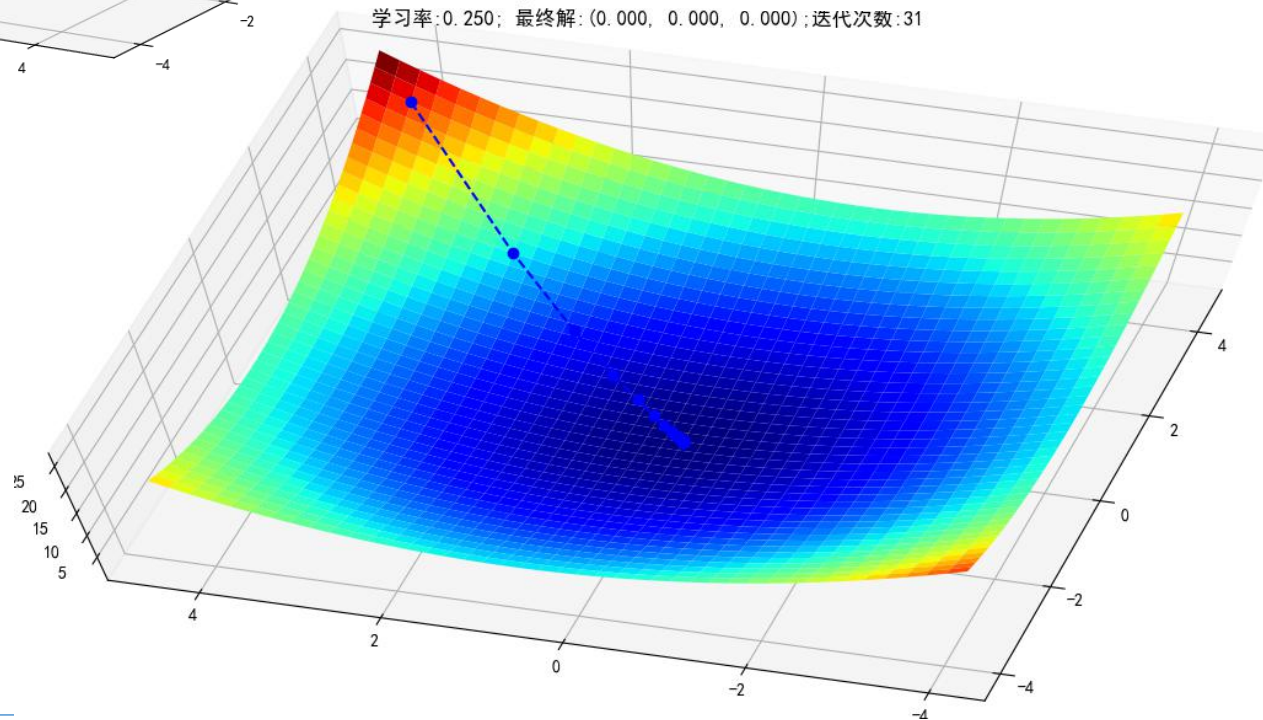
函数 $y = 0.6 * (\theta_1 + \theta_2)^2 - \theta_1 * \theta_2$;
学习率: 0.250; 最终解: (0.000, 0.000, 0.000); 迭代次数: 31



函数 $y = 0.6 * (\theta_1 + \theta_2)^2 - \theta_1 * \theta_2$;
学习率: 0.100; 最终解: (0.000, 0.000, 0.000); 迭代次数: 22

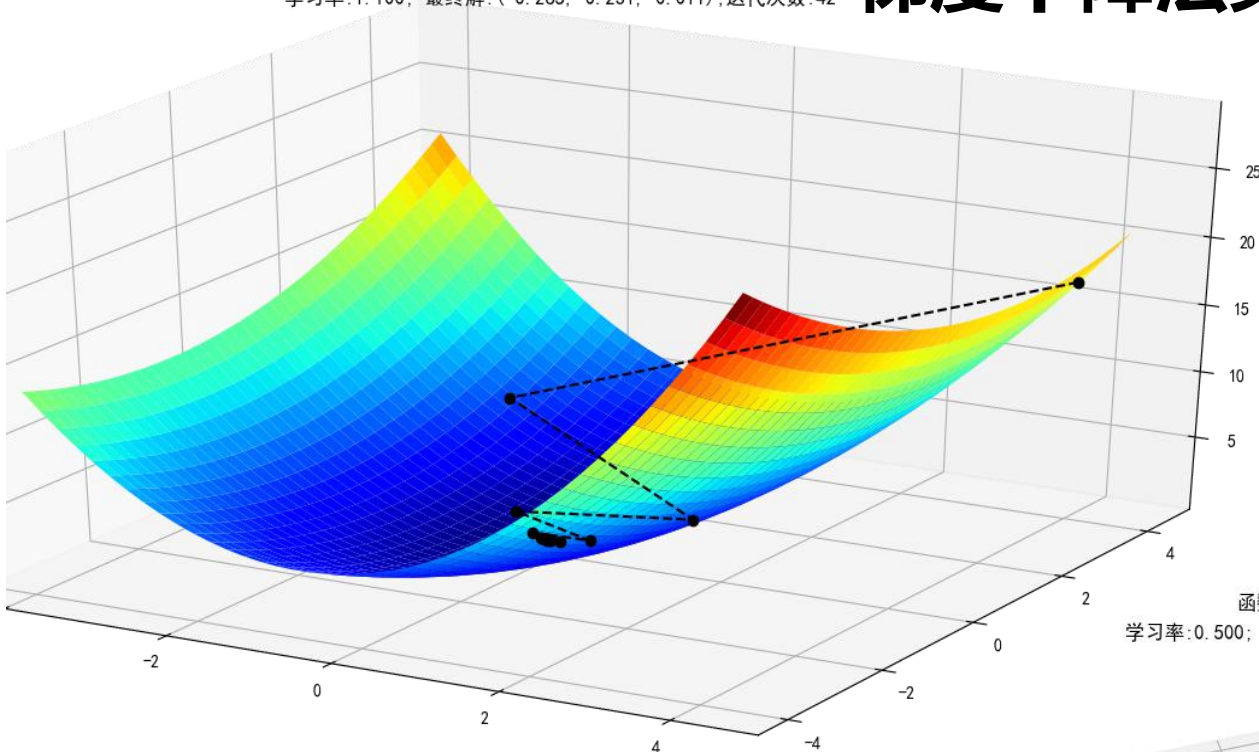


函数 $y = 0.6 * (\theta_1 + \theta_2)^2 -$
学习率: 0.250; 最终解: (0.000, 0.000, 0.000); 迭代次数: 31

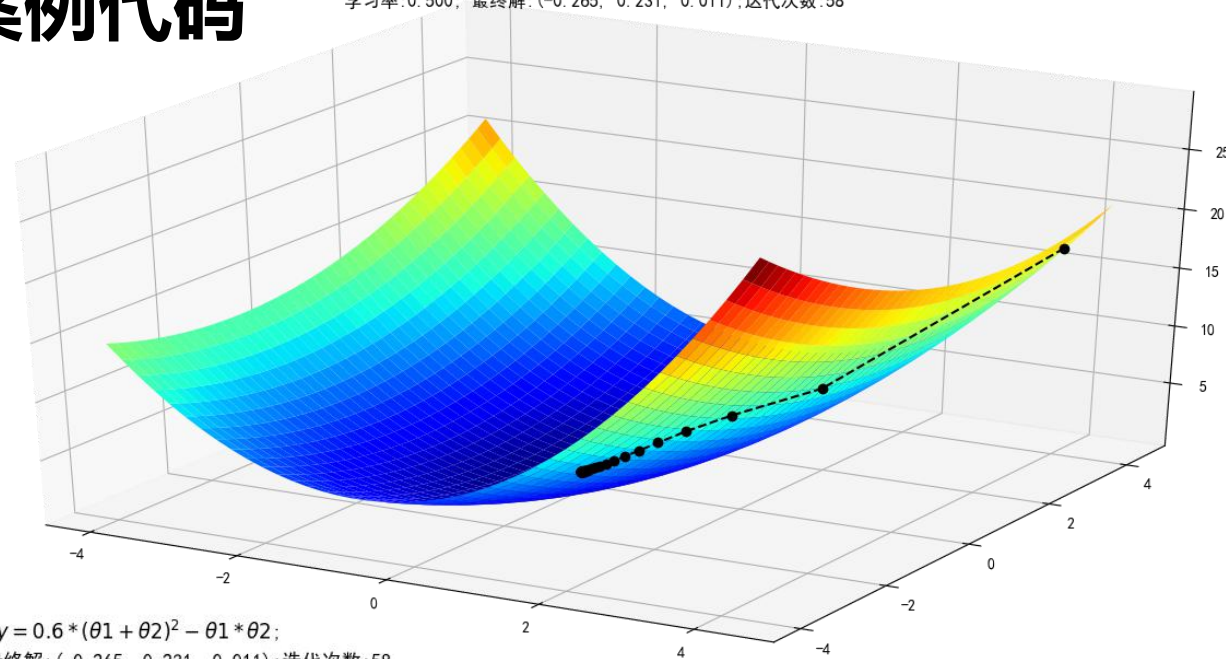


梯度下降法案例代码

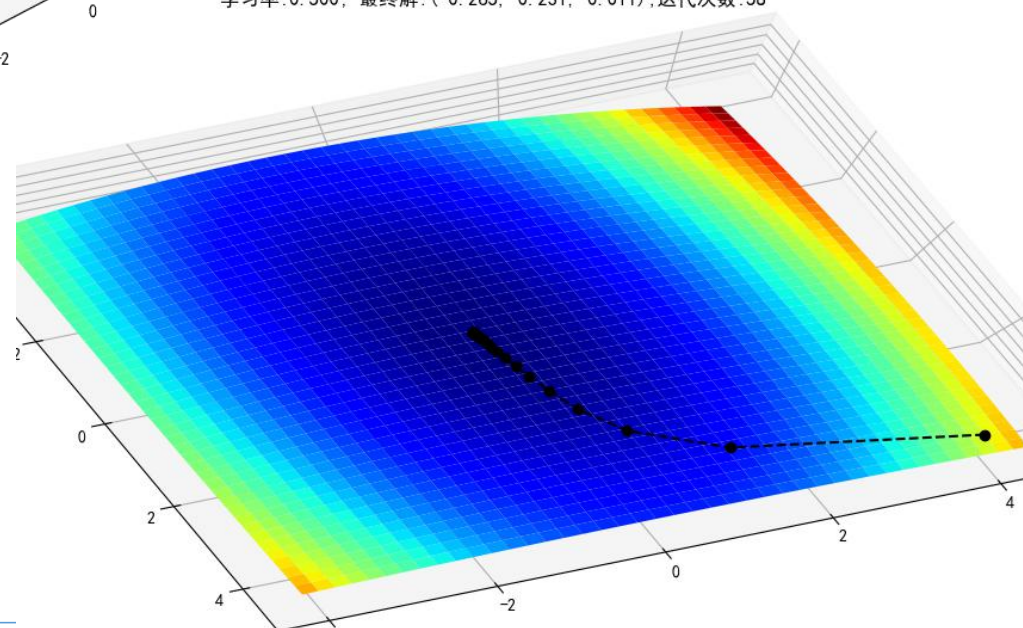
函数 $y = 0.6 * (\theta_1 + \theta_2)^2 - \theta_1 * \theta_2$;
学习率: 1.100; 最终解: (-0.265, 0.231, 0.011); 迭代次数: 42



函数 $y = 0.6 * (\theta_1 + \theta_2)^2 - \theta_1 * \theta_2$;
学习率: 0.500; 最终解: (-0.265, 0.231, 0.011); 迭代次数: 58



函数 $y = 0.6 * (\theta_1 + \theta_2)^2 - \theta_1 * \theta_2$;
学习率: 0.500; 最终解: (-0.265, 0.231, 0.011); 迭代次数: 58



梯度下降法

- 由于梯度下降法中负梯度方向作为变量的变化方向，所以有可能导致最终求解的值是局部最优解，所以在使用梯度下降的时候，一般需要进行一些调优策略：
 - ◆ 学习率的选择：学习率过大，表示每次迭代更新的时候变化比较大，有可能会跳过最优解；学习率过小，表示每次迭代更新的时候变化比较小，就会导致迭代速度过慢，很长时间都不能结束；
 - ◆ 算法初始参数值的选择：初始值不同，最终获得的最小值也有可能不同，因为梯度下降法求解的是局部最优解，所以一般情况下，选择多次不同初始值运行算法，并最终返回损失函数最小情况下的结果值；
 - ◆ 标准化：由于样本不同特征的取值范围不同，可能会导致在各个不同参数上迭代速度不同，为了减少特征取值的影响，可以将特征进行标准化操作。

梯度下降法

- 批量梯度下降法(Batch Gradient Descent, BGD)：使用所有样本在当前点的梯度值来对变量参数进行更新操作。

$$\theta^{k+1} = \theta^k - \alpha \sum_{i=1}^m \nabla f_{\theta^k}(x^i)$$

- 随机梯度下降法(Stochastic Gradient Descent, SGD)：在更新变量参数的时候，选取一个样本的梯度值来更新参数。

$$\theta^{k+1} = \theta^k - \alpha \nabla f_{\theta^k}(x^i)$$

- 小批量梯度下降法(Mini-batch Gradient Descent, MBGD)：结合BGD和SGD的特性，从原始数据中，每次选择n个样本来更新参数值，一般n选择10。

$$\theta^{k+1} = \theta^k - \alpha \sum_{i=t}^{t+n-1} \nabla f_{\theta^k}(x^i)$$

梯度下降法

■ BGD、SGD、MBGD的区别：

- ◆ 当样本量为 m 的时候，每次迭代BGD算法中对于参数值更新一次，SGD算法中对于参数值更新 m 次，MBGD算法中对于参数值更新 m/n 次，相对来讲SGD算法的更新速度最快；
- ◆ SGD算法中对于每个样本都需要更新参数值，当样本值不太正常的时候，就有可能导致本次的参数更新会产生相反的影响，也就是说SGD算法的结果并不是完全收敛的，而是在收敛结果处波动的；
- ◆ SGD算法是每个样本都更新一次参数值，所以SGD算法特别适合样本数据量大的情况以及在线机器学习(Online ML)。

Python科学计算库回顾

- Python科学计算库主要是为机器学习提供了一些便捷、封装好的API，那么在实际工作中，主要是将其应用在机器学习的特征工程阶段，主要涉及到的库有以下几个：
 - ◆ **NumPy**-数学计算基础库：N维数组、线性代数计算、傅立叶变换、随机数等；
 - ◆ **SciPy**-数值计算库：线性代数、拟合与优化、插值、数值积分、稀疏矩阵、图像处理、统计等；
 - ◆ **Pandas**-数据分析库：数据导入、整理、处理、分析等；
 - ◆ **Matplotlib**-绘图库：绘制二维图形和图表。

Python科学计算库回顾

- 官网：<https://scipy.org/>



NumPy
Base N-dimensional
array package



SciPy library
Fundamental library
for scientific
computing



Matplotlib
Comprehensive 2D
Plotting



IPython
Enhanced Interactive
Console



Sympy
Symbolic mathematics



pandas
Data structures &
analysis

- 参考文档：

- ◆ <http://python.usyiyi.cn/>
- ◆ <https://docs.scipy.org/doc/>
- ◆ <http://pandas.pydata.org/pandas-docs/stable/index.html>



THANK YOU

上海育创网络科技有限公司