

# 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，小象学院和主讲老师拥有完全知识产权的权利；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。任何其他人或机构不得盗版、复制、仿造其中的创意及内容，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：小象

■ 新浪微博：ChinaHadoop



# Python基础

---



小象学院  
ChinaHadoop.cn

邹博

# 本次说明

---

- 本PPT后面仅列举使用Python库的效果截图，详细内容请参考该PPT的配套代码。

# Python库

---

- Pip
  - 安装Python包的推荐工具: <https://pypi.python.org/pypi/pip>
- Numpy
  - 为Python提供快速的多维数组处理能力
- Pandas: Python Data Analysis Library
  - 在Numpy基础上提供了更多的数据读写工具
- Scipy
  - 在NumPy基础上添加了众多科学计算工具包
- Matplotlib
  - Python丰富的绘图库
- 官网:
  - Numpy/Scipy: <http://www.scipy.org>
  - Pandas: <http://pandas.pydata.org/>
  - Matplotlib: <http://www.matplotlib.org>

Country	Year	Value
Algeria	2000	0.0000
Algeria	2001	0.0000
Algeria	2002	0.0000
Algeria	2003	0.0000
Algeria	2004	0.0000
Algeria	2005	0.0000
Algeria	2006	0.0000
Algeria	2007	0.0000
Algeria	2008	0.0000
Algeria	2009	0.0000
Algeria	2010	0.0000
Algeria	2011	0.0000
Algeria	2012	0.0000
Algeria	2013	0.0000
Algeria	2014	0.0000
Algeria	2015	0.0000
Algeria	2016	0.0000
Algeria	2017	0.0000
Algeria	2018	0.0000
Algeria	2019	0.0000
Algeria	2020	0.0000
Algeria	2021	0.0000
Algeria	2022	0.0000
Algeria	2023	0.0000
Algeria	2024	0.0000
Algeria	2025	0.0000
Algeria	2026	0.0000
Algeria	2027	0.0000
Algeria	2028	0.0000
Algeria	2029	0.0000
Algeria	2030	0.0000
Algeria	2031	0.0000
Algeria	2032	0.0000
Algeria	2033	0.0000
Algeria	2034	0.0000
Algeria	2035	0.0000
Algeria	2036	0.0000
Algeria	2037	0.0000
Algeria	2038	0.0000
Algeria	2039	0.0000
Algeria	2040	0.0000
Algeria	2041	0.0000
Algeria	2042	0.0000
Algeria	2043	0.0000
Algeria	2044	0.0000
Algeria	2045	0.0000
Algeria	2046	0.0000
Algeria	2047	0.0000
Algeria	2048	0.0000
Algeria	2049	0.0000
Algeria	2050	0.0000
Algeria	2051	0.0000
Algeria	2052	0.0000
Algeria	2053	0.0000
Algeria	2054	0.0000
Algeria	2055	0.0000
Algeria	2056	0.0000
Algeria	2057	0.0000
Algeria	2058	0.0000
Algeria	2059	0.0000
Algeria	2060	0.0000
Algeria	2061	0.0000
Algeria	2062	0.0000
Algeria	2063	0.0000
Algeria	2064	0.0000
Algeria	2065	0.0000
Algeria	2066	0.0000
Algeria	2067	0.0000
Algeria	2068	0.0000
Algeria	2069	0.0000
Algeria	2070	0.0000
Algeria	2071	0.0000
Algeria	2072	0.0000
Algeria	2073	0.0000
Algeria	2074	0.0000
Algeria	2075	0.0000
Algeria	2076	0.0000
Algeria	2077	0.0000
Algeria	2078	0.0000
Algeria	2079	0.0000
Algeria	2080	0.0000
Algeria	2081	0.0000
Algeria	2082	0.0000
Algeria	2083	0.0000
Algeria	2084	0.0000
Algeria	2085	0.0000
Algeria	2086	0.0000
Algeria	2087	0.0000
Algeria	2088	0.0000
Algeria	2089	0.0000
Algeria	2090	0.0000
Algeria	2091	0.0000
Algeria	2092	0.0000
Algeria	2093	0.0000
Algeria	2094	0.0000
Algeria	2095	0.0000
Algeria	2096	0.0000
Algeria	2097	0.0000
Algeria	2098	0.0000
Algeria	2099	0.0000
Algeria	2100	0.0000
Algeria	2101	0.0000
Algeria	2102	0.0000
Algeria	2103	0.0000
Algeria	2104	0.0000
Algeria	2105	0.0000
Algeria	2106	

```
Collecting pandas
  Downloading pandas-0.19.1-cp27-cp27m-win32.whl (6.7MB)
    100% |████████████████████████████████████████████████████████████████████████████████| 6.7MB 81kB/s
Collecting python-dateutil (from pandas)
  Downloading python_dateutil-2.6.0-py2.py3-none-any.whl (194kB)
    100% |████████████████████████████████████████████████████████████████████████████████| 194kB 103kB/s
Collecting numpy>=1.7.0 (from pandas)
  Downloading numpy-1.11.2-cp27-none-win32.whl (6.5MB)
    100% |████████████████████████████████████████████████████████████████████████████████| 6.5MB 103kB/s
Collecting pytz>=2011k (from pandas)
  Downloading pytz-2016.7-py2.py3-none-any.whl (480kB)
    100% |████████████████████████████████████████████████████████████████████████████████| 481kB 244kB/s
Requirement already up-to-date: six>=1.5 in c:\python27\lib\site-packages (from python-dateutil->pandas)
Installing collected packages: python-dateutil, numpy, pytz, pandas
  Found existing installation: python-dateutil 2.5.2
    Uninstalling python-dateutil-2.5.2:
      Successfully uninstalled python-dateutil-2.5.2
  Found existing installation: numpy 1.11.1+mk1
    Uninstalling numpy-1.11.1+mk1:
      Successfully uninstalled numpy-1.11.1+mk1
  Found existing installation: pytz 2016.3
    Uninstalling pytz-2016.3:
      Successfully uninstalled pytz-2016.3
  Found existing installation: pandas 0.18.1
    Uninstalling pandas-0.18.1:
      Successfully uninstalled pandas-0.18.1
Successfully installed numpy-1.11.2 pandas-0.19.1 python-dateutil-2.6.0 pytz-2016.7
```

## Collecting pandas

Downloading pandas-0.19.1-cp27-cp27m-win32.whl (6.7MB)

100% ██████████ 6.7MB 81kB/s

## Collecting python-dateutil (from pandas)

Downloading python-dateutil-2.6.0-py2.py3-none-any.whl (194kB)

100% ██████████ 194kB 103kB/s

Collecting numpy&gt;=1.7.0 (from pandas)

Downloading numpy-1.11.2-cp27-none-win32.whl (6.5MB)

100% 6.5MB 103kB/s

Collecting pytz>=2011k (from pandas)

Downloading pytz-2016.7-py2.py3-none-any.whl (480kB)

100% ██████████ 481kB 244kB/s

Requirement already up-to-date: six>=1.5 in c:\python27\lib\site-packages (from python-dateutil->pandas)

```
Installing collected packages: python-dateutil, numpy, pytz, pandas
```

```
Found existing installation: python-dateutil 2.5.2
```

### Uninstalling python-dateutil-2.5.2:

Successfully uninstalled python-dateutil-2.5.2

```
Found existing installation: numpy 1.11.1+mkl
```

Uninstalling numpy-1.11.1+mk1:

Successfully uninstalled numpy-1.11.1+mk1

```
Found existing installation: pvtz 2016.3
```

### Uninstalling pytz-2016.3:

Successfully uninstalled pytz-2016.3

```
Found existing installation: pandas 0.18.1
```

Uninstalling pandas-0.18.1:

Successfully uninstalled pandas-0.18.1

Successfully installed numpy-1.11.2 pandas-0.19.1 python-dateutil-2.6.0 pytz-2016.7

# 数据生成

---

□  $a = np.arange(0, 60, 10).reshape((-1, 1)) + np.arange(6)$

□  $A =$

```
[[ 0  1  2  3  4  5]
 [10 11 12 13 14 15]
 [20 21 22 23 24 25]
 [30 31 32 33 34 35]
 [40 41 42 43 44 45]
 [50 51 52 53 54 55]]
```

# Taylor展式的应用

```
def calc_e_small(x):
    n = 10
    f = np.arange(1, n+1).cumprod()
    b = np.array([x]*n).cumprod()
    return np.sum(b / f) + 1

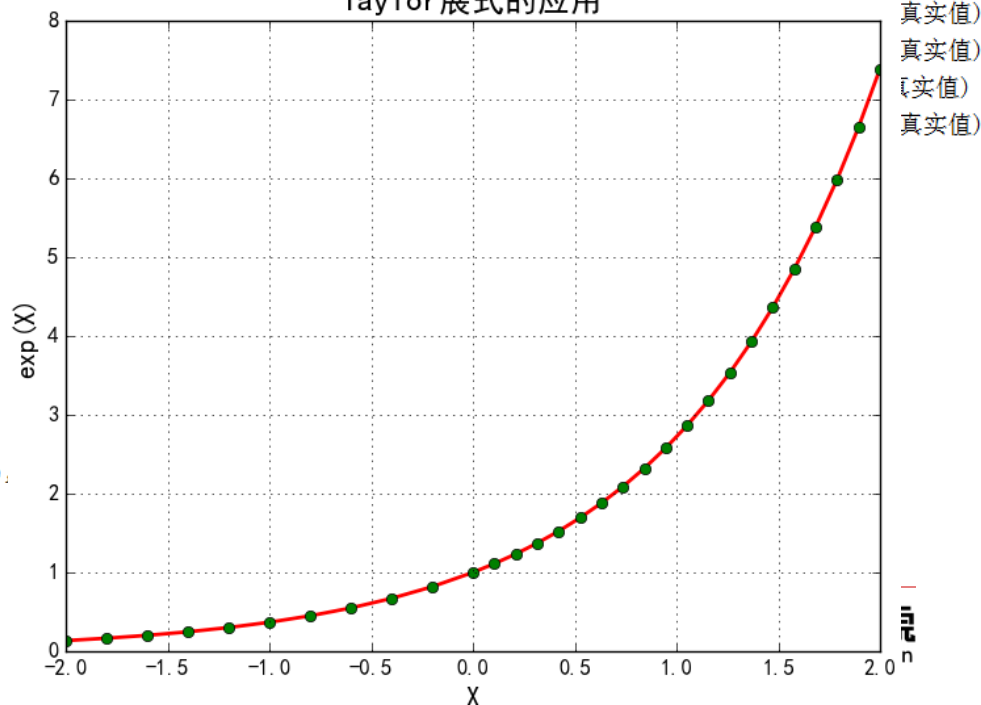
def calc_e(x):
    reverse = False
    if x < 0: # 处理负数
        x = -x
        reverse = True
    ln2 = 0.69314718055994530941723212145818
    c = x / ln2
    a = int(c+0.5)
    b = x - a*ln2
    y = (2 ** a) * calc_e_small(b)
    if reverse:
        return 1/y
    return y

if __name__ == "__main__":
    t1 = np.linspace(-2, 0, 10, endpoint=False)
    t2 = np.linspace(0, 2, 20)
    t = np.concatenate((t1, t2))
    print t # 横轴数据
    y = np.empty_like(t)
    for i, x in enumerate(t):
        y[i] = calc_e(x)
        print 'e^', x, ' = ', y[i], '(近似值)\t', math.exp(x),
        # print '误差: ', y[i] - math.exp(x)

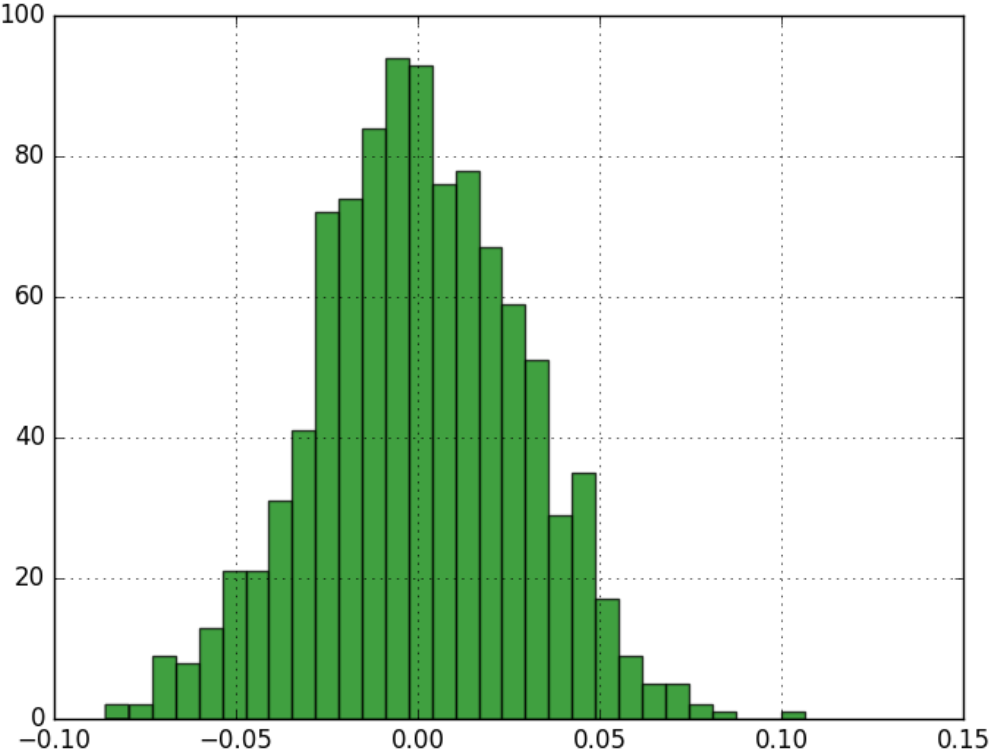
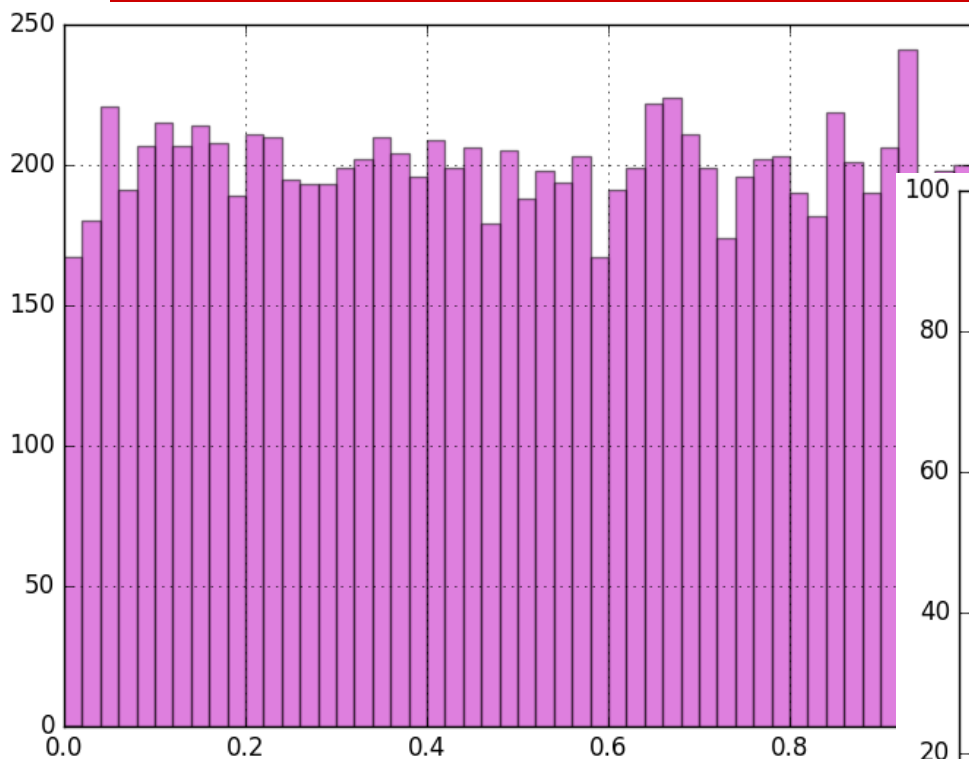
    mpl.rcParams['font.sans-serif'] = [u'SimHei']
    mpl.rcParams['axes.unicode_minus'] = False
    plt.plot(t, y, 'r-', t, y, 'go', linewidth=2)
    plt.title(u'Taylor展式的应用', fontsize=18)
    plt.xlabel('X', fontsize=15)
    plt.ylabel('exp(X)', fontsize=15)
    plt.grid(True)
    plt.show()
```

$e^{-0.8}$  = 0.449328964117 (近似值) 0.449328964117 (真实值)  
 $e^{-0.6}$  = 0.548811636094 (近似值) 0.548811636094 (真实值)  
 $e^{-0.4}$  = 0.670320046036 (近似值) 0.670320046036 (真实值)  
 $e^{-0.2}$  = 0.818730753078 (近似值) 0.818730753078 (真实值)  
 $e^{0.0}$  = 1.0 (近似值) 1.0 (真实值)  
 $e^{0.105263157895}$  = 1.11100294108 (近似值) 1.11100294108 (真实值)  
 $e^{0.210526315789}$  = 1.2343275351 (近似值) 1.2343275351 (真实值)  
 $e^{0.315789473684}$  = 1.37134152176 (近似值) 1.37134152176 (真实值)  
 $e^{0.421052631579}$  = 1.5235644639 (近似值) 1.5235644639 (真实值)  
 $e^{0.526315789474}$  = 1.69268460033 (近似值) 1.69268460033 (真实值)  
 $e^{0.631578947368}$  = 1.88057756929 (近似值) 1.88057756929 (真实值)  
 $e^{0.736842105263}$  = 2.08932721042 (近似值) 2.08932721042 (真实值)  
 $e^{0.842105263158}$  = 2.32124867566 (近似值) 2.32124867566 (真实值)  
 $e^{0.947368421053}$  = 2.57891410565 (近似值) 2.57891410565 (真实值)  
 $e^{1.05263157895}$  = 2.86518115618 (近似值) 2.86518115618 (真实值)  
 $e^{1.15789473684}$  = 3.18322469126 (近似值) 3.18322469126 (真实值)  
 $e^{1.26315789474}$  = 3.53657199412 (近似值) 3.53657199412 (真实值)  
 $e^{1.36842105263}$  = 3.92914188683 (近似值) 3.92914188683 (真实值)  
 $e^{1.47368421053}$  = 4.3652881922 (近似值) 4.3652881922 (真实值)

Taylor展式的应用

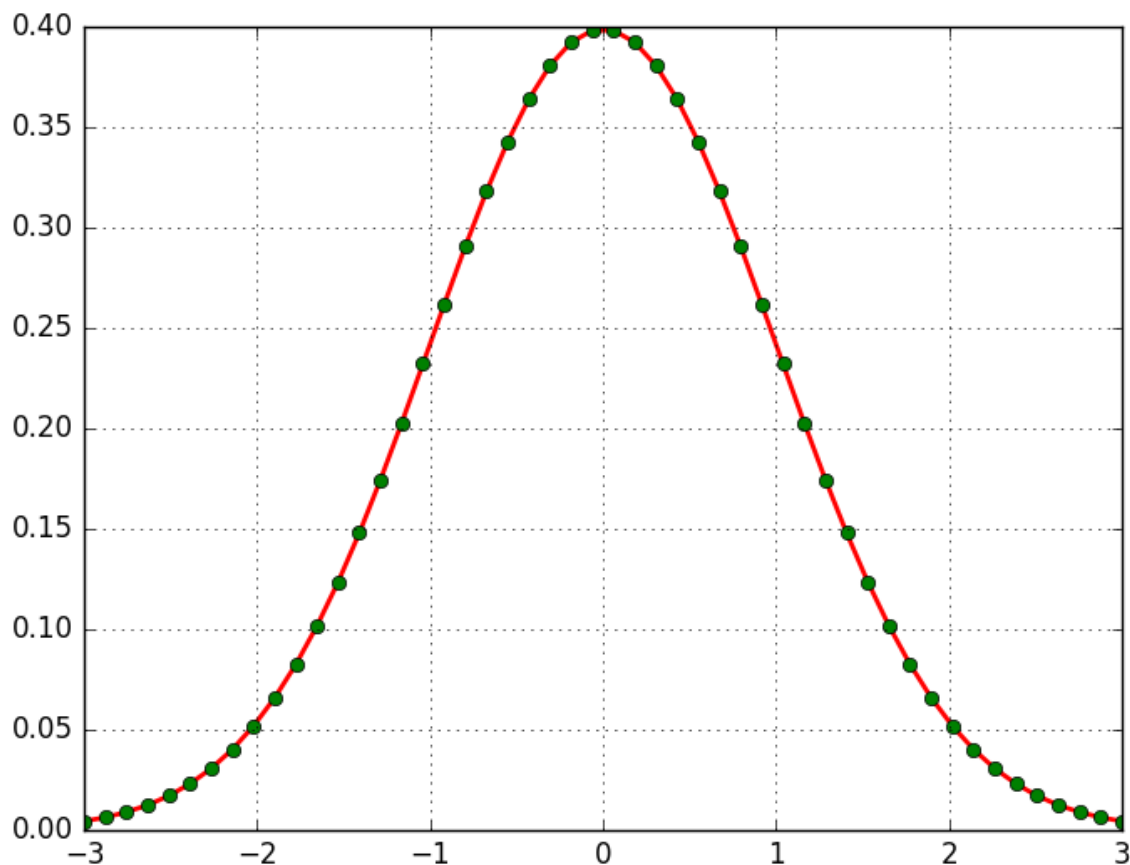


# 验证中心极限定理

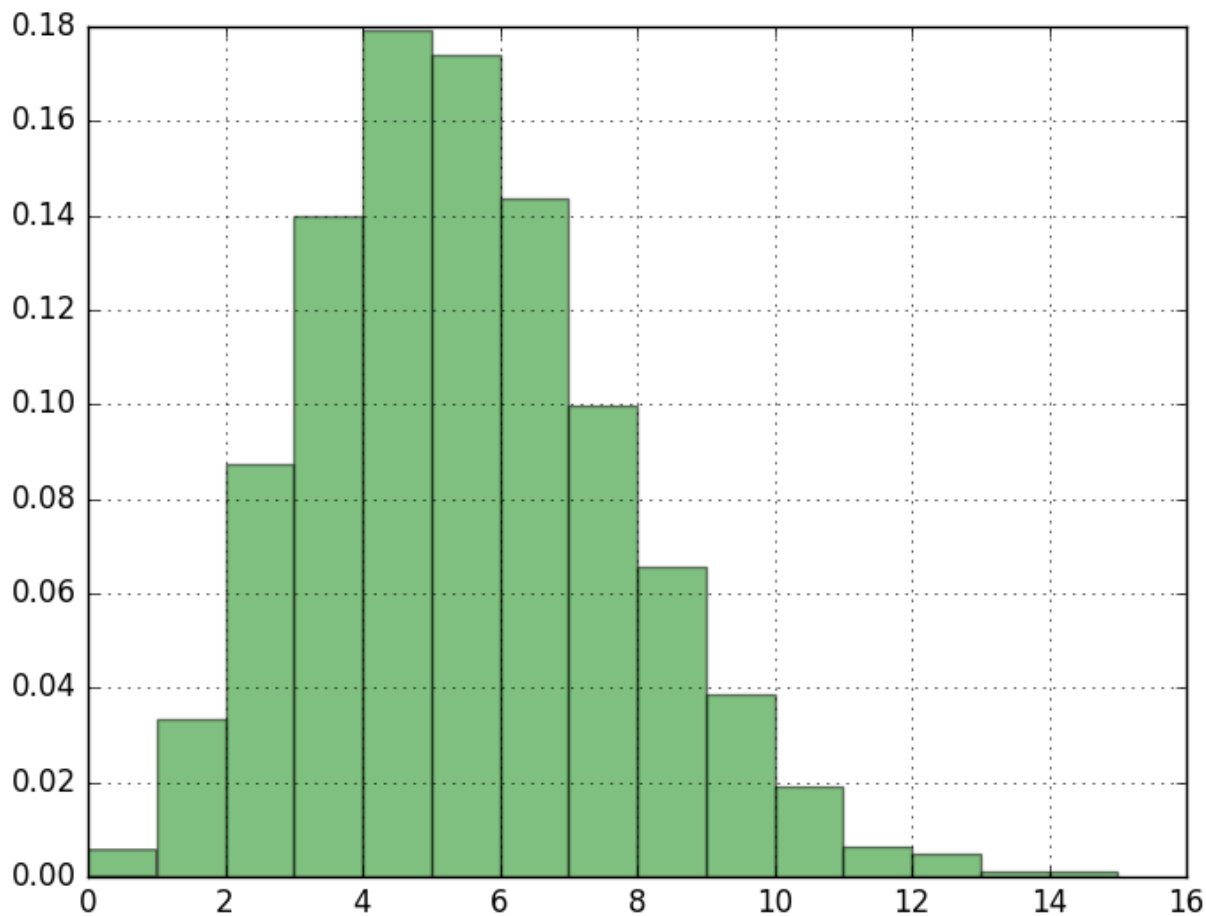




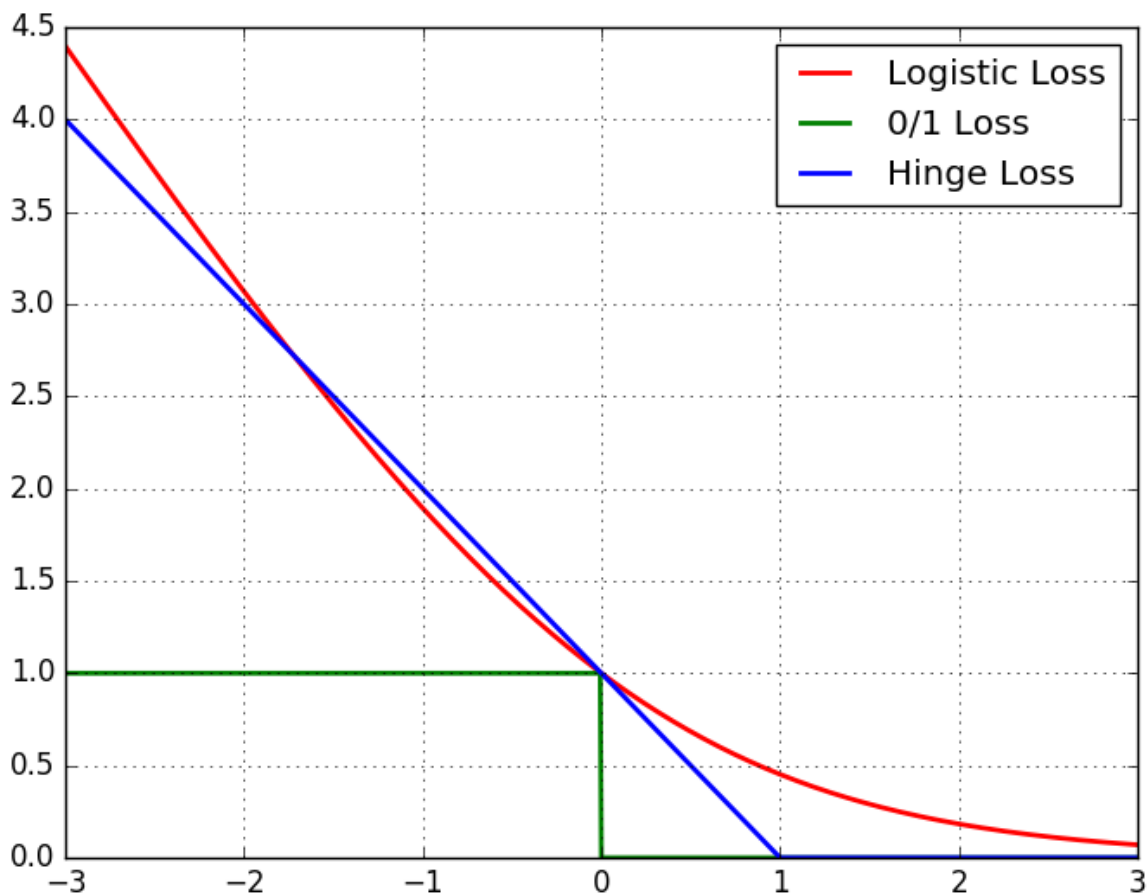
# 正态分布的概率密度函数



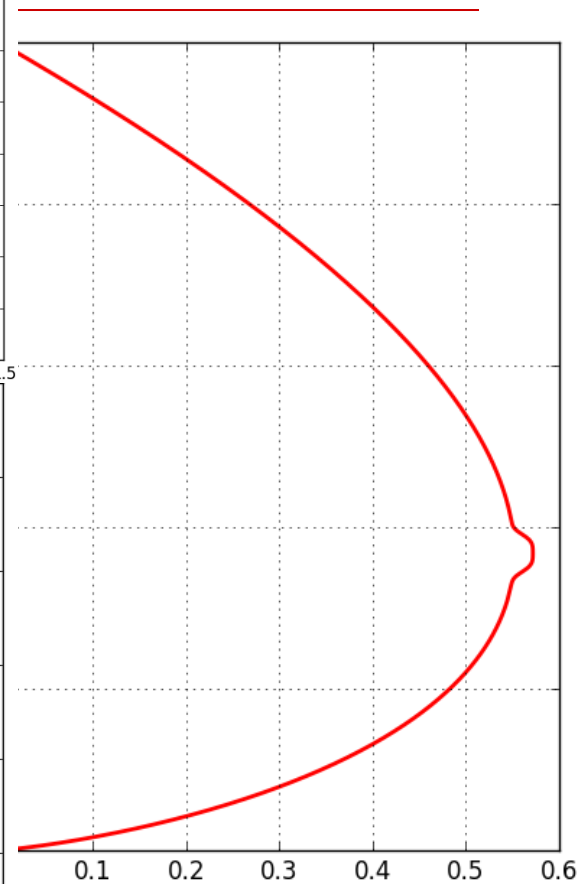
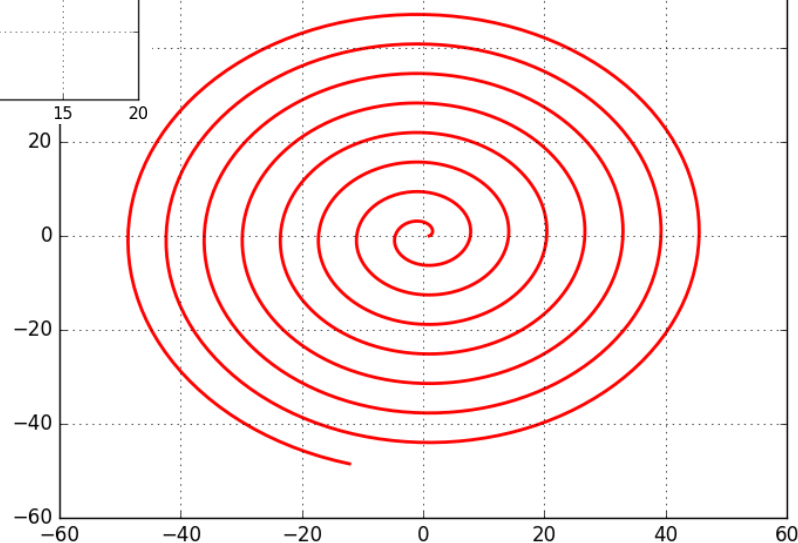
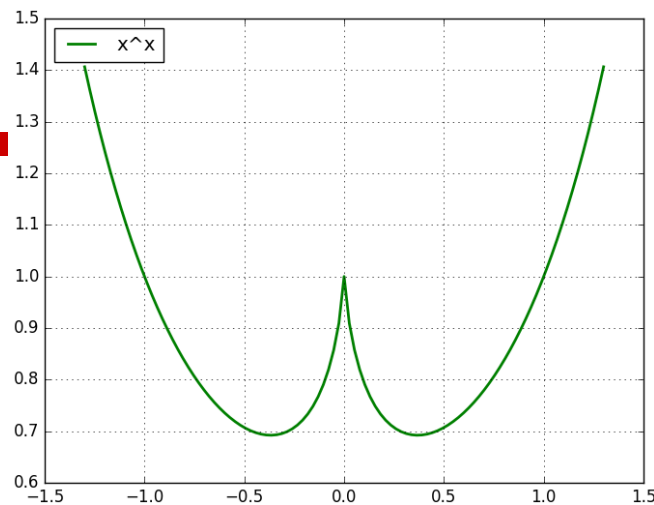
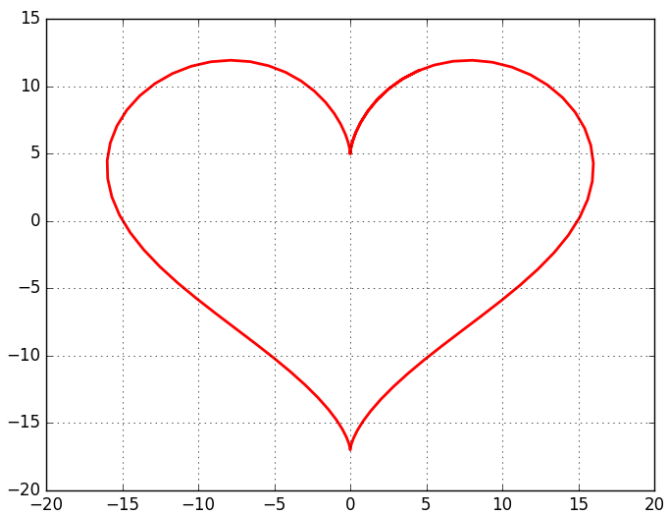
# Poisson分布的概率质量函数



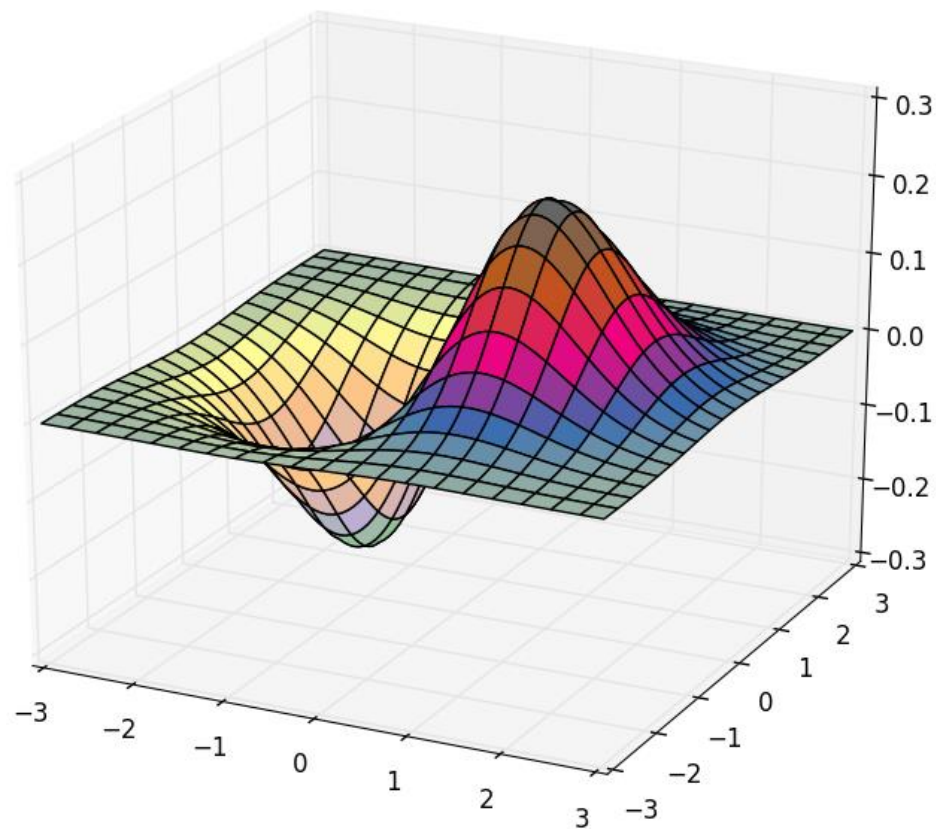
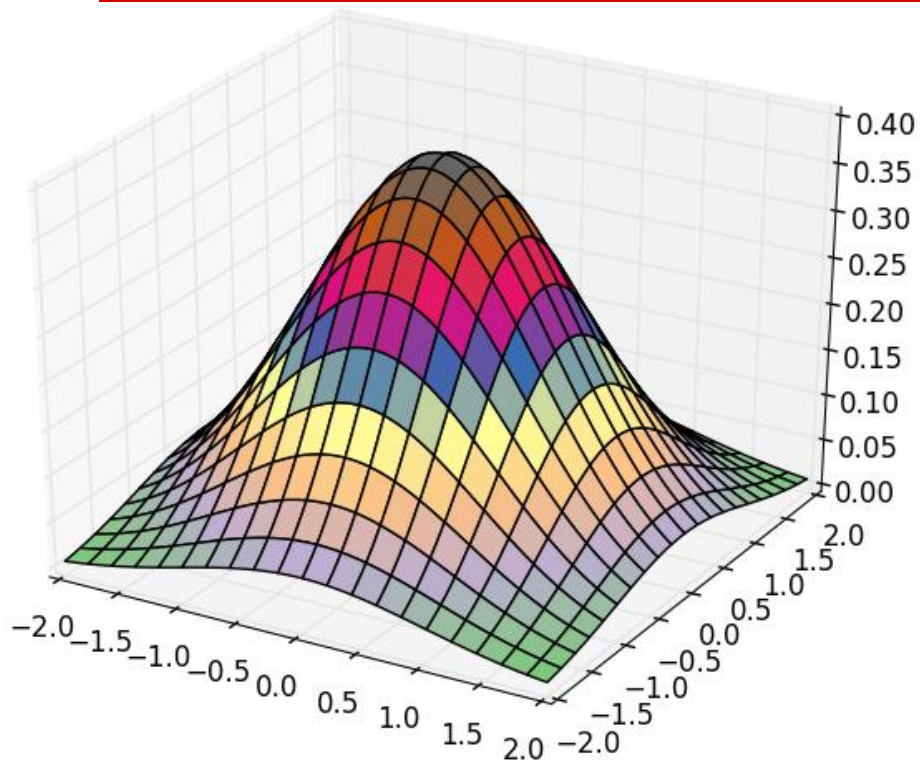
# 机器学习中的损失函数



# 各种2D曲线



# 3D



# 类/继承类

```
class People:
    def __init__(self, n, a, s):
        self.name = n
        self.age = a
        self.__score = s
        self.print_people()
        # self.__print_people() # 私有函数的作用

    def print_people(self):
        str = u'%s的年龄: %d, 成绩为: %.2f' % (self.name, self.age, self.__score)
        print str

    __print_people = print_people

class Student(People):
    def __init__(self, n, a, w):
        People.__init__(self, n, a, w)
        self.name = 'Student ' + self.name

    def print_people(self):
        str = u'%s的年龄: %d' % (self.name, self.age)
        print str

def func(p):
    p.age = 11

if __name__ == '__main__':
    p = People('Tom', 10, 3.14159)
    func(p) # p传入的是引用类型
    p.print_people()

    # 注意分析下面语句的打印结果, 是否觉得有些“怪异”?
    j = Student('Jerry', 12, 2.71828)

    # 成员函数
    j.print_people()
    People.print_people(j)
```

Tom的年龄: 10, 成绩为: 3.14

Tom的年龄: 11, 成绩为: 3.14

Jerry的年龄: 12

Tom的年龄: 11, 成绩为: 3.14

Student Jerry的年龄: 12

Tom的年龄: 11, 成绩为: 3.14

Student Jerry的年龄: 12, 成绩为: 2.72

# 统计量

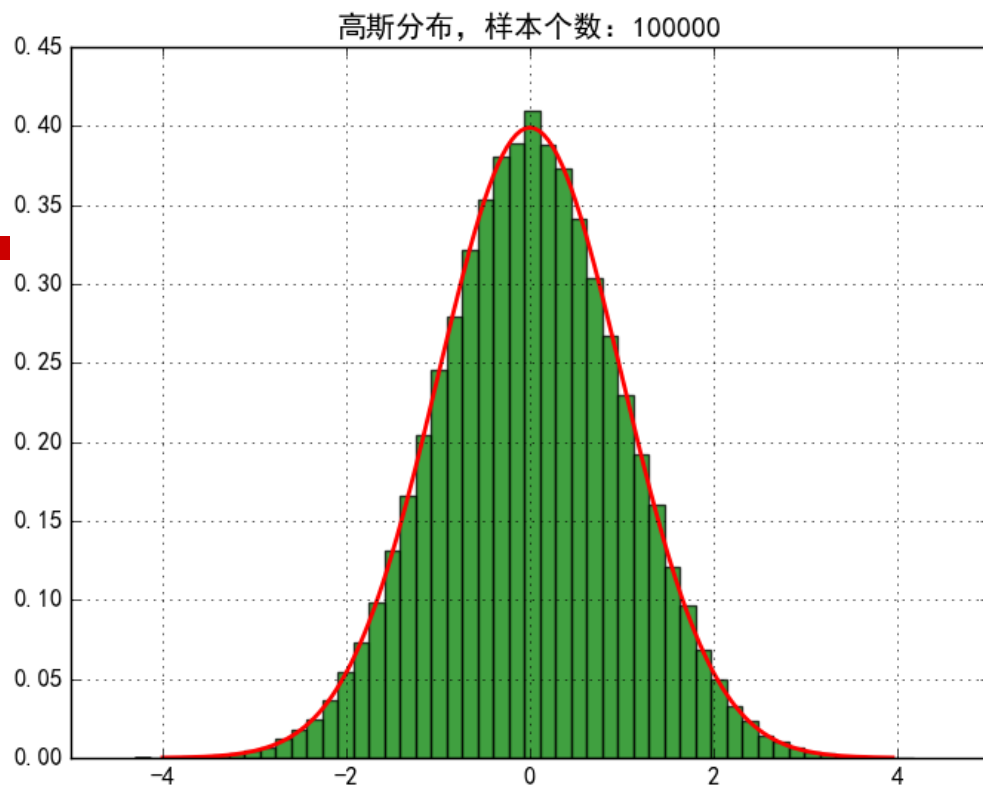
```
def calc_statistics(x):
    n = x.shape[0] # 样本个数

    # 手动计算
    m = 0
    m2 = 0
    m3 = 0
    m4 = 0
    for t in x:
        m += t
        m2 += t*t
        m3 += t**3
        m4 += t**4

    m /= n
    m2 /= n
    m3 /= n
    m4 /= n

    mu = m
    sigma = np.sqrt(m2 - mu*mu)
    skew = (m3 - 3*mu*m2 + 2*mu**3) / sigma**3
    kurtosis = (m4 - 4*mu*m3 + 6*mu*mu*m2 - 4*mu**3*mu + mu**4) / sigma**4 - 3
    print '手动计算均值、标准差、偏度、峰度: ', mu, sigma, skew, kurtosis

    # 使用系统函数验证
    mu = np.mean(x, axis=0)
    sigma = np.std(x, axis=0)
    skew = stats.skew(x)
    kurtosis = stats.kurtosis(x)
    return mu, sigma, skew, kurtosis
```

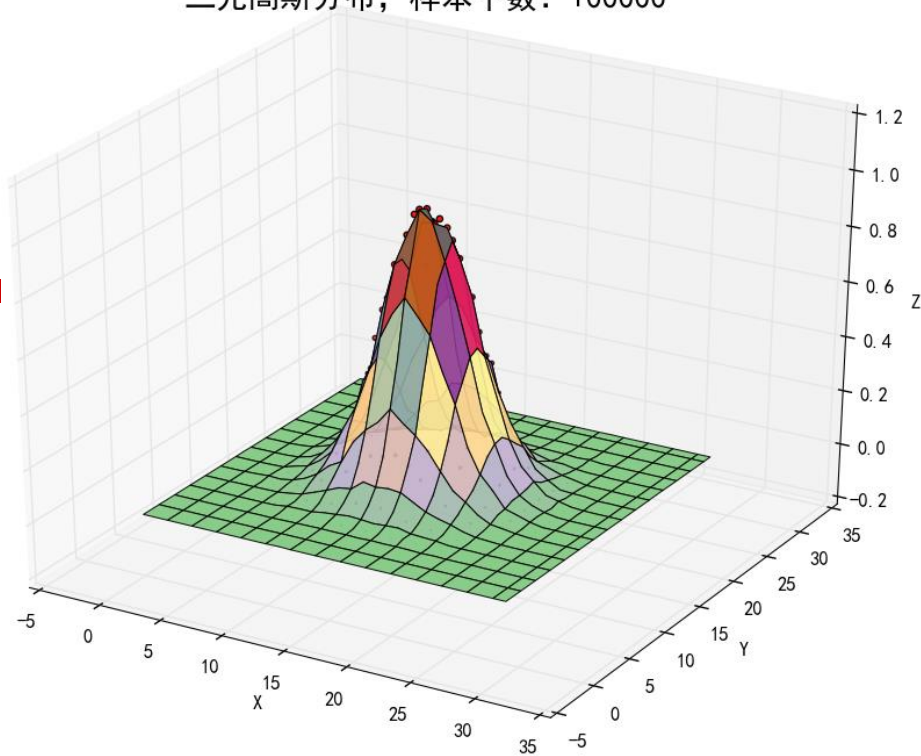


手动计算均值、标准差、偏度、峰度： -0.00232018730484 1.00220229337 0.0070687774347 0.0174102810253

函数库计算均值、标准差、偏度、峰度： -0.00232018730484 1.00220229337 0.0070687774347 0.0174102810253

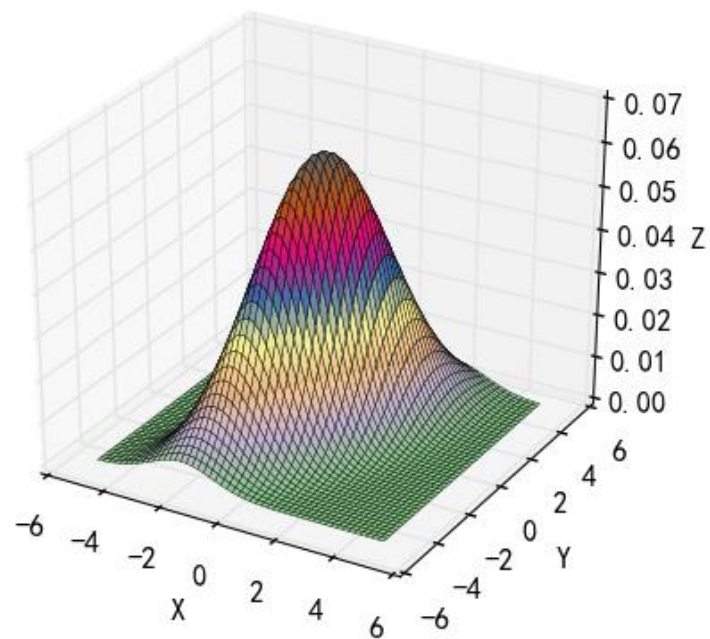
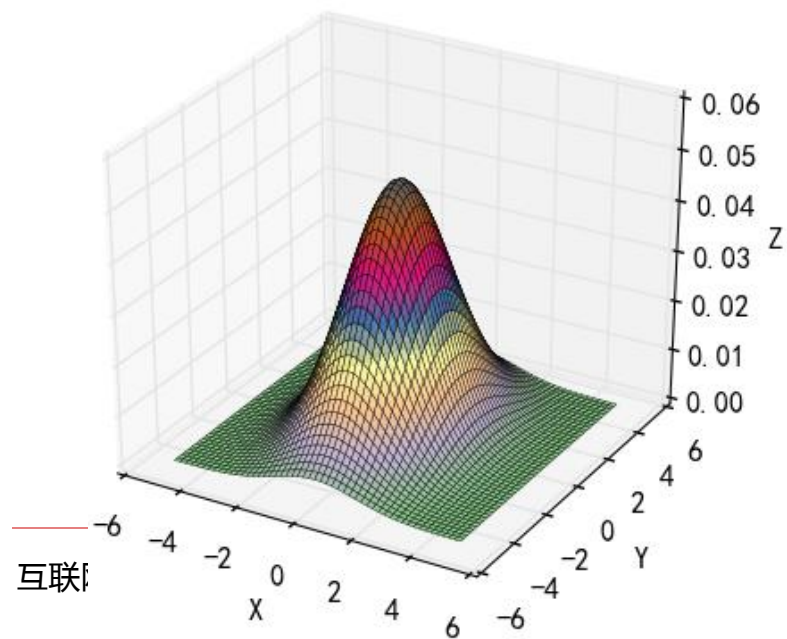
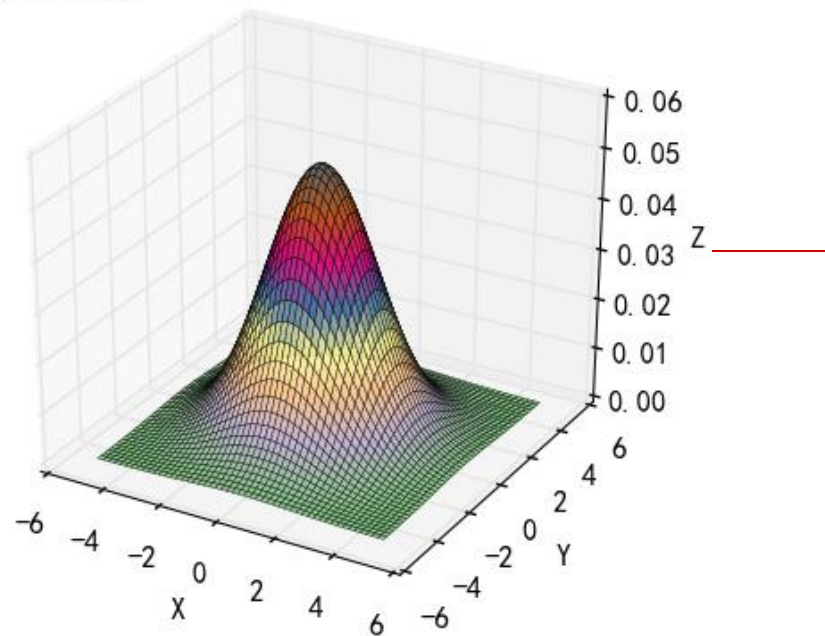
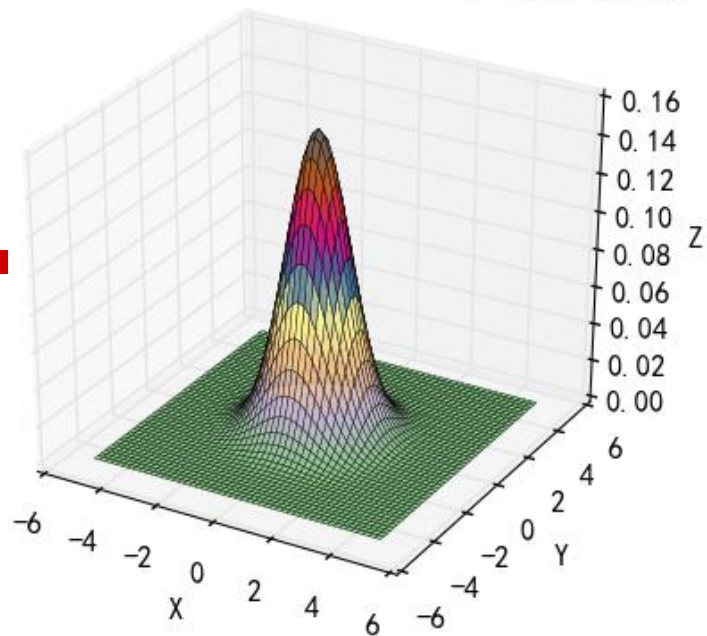
# 二元高斯分布

```
d = np.random.randn(100000, 2)
mu, sigma, skew, kurtosis = calc_statistics(d)
print '函数库计算均值、标准差、偏度、峰度: ', mu,
# 二维图像
N = 30
density, edges = np.histogramdd(d, bins=[N, N])
print '样本总数: ', np.sum(density)
density /= density.max()
x = y = np.arange(N)
t = np.meshgrid(x, y)
fig = plt.figure(facecolor='w')
ax = fig.add_subplot(111, projection='3d')
ax.scatter(t[0], t[1], density, c='r', s=15*density, marker='o', depthshade=True)
ax.plot_surface(t[0], t[1], density, cmap=cm.Accent, rstride=2, cstride=2, alpha=0.9, lw=0.75)
ax.set_xlabel(u'X')
ax.set_ylabel(u'Y')
ax.set_zlabel(u'Z')
plt.title(u'二元高斯分布，样本个数: %d' % d.shape[0], fontsize=20)
plt.tight_layout(0.1)
plt.show()
```

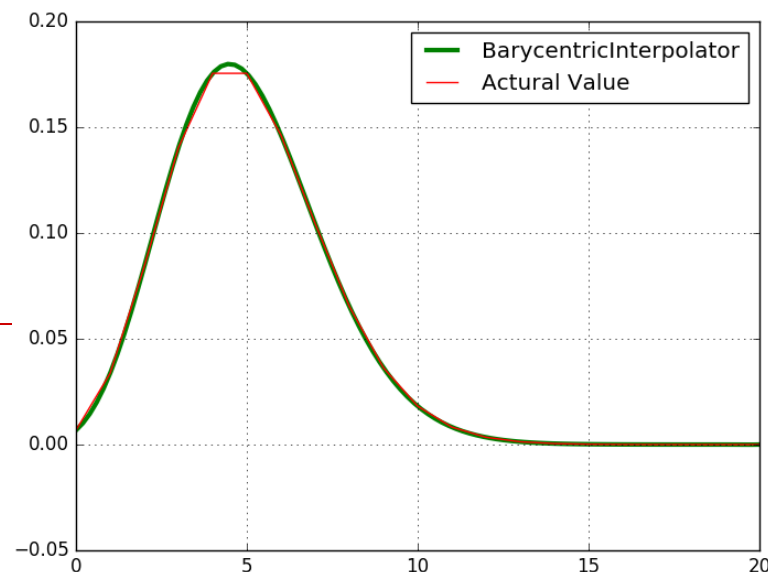




## 二元高斯分布方差比较



# 重心插值



□ 给定实数对  $\{(x_j, y_j), j = 0, 1, \dots, n\}$

■  $x_j$  互不相同。

□ 对于给定的  $n+1$  个权值  $\{u_j \neq 0, j = 0, 1, \dots, n\}$

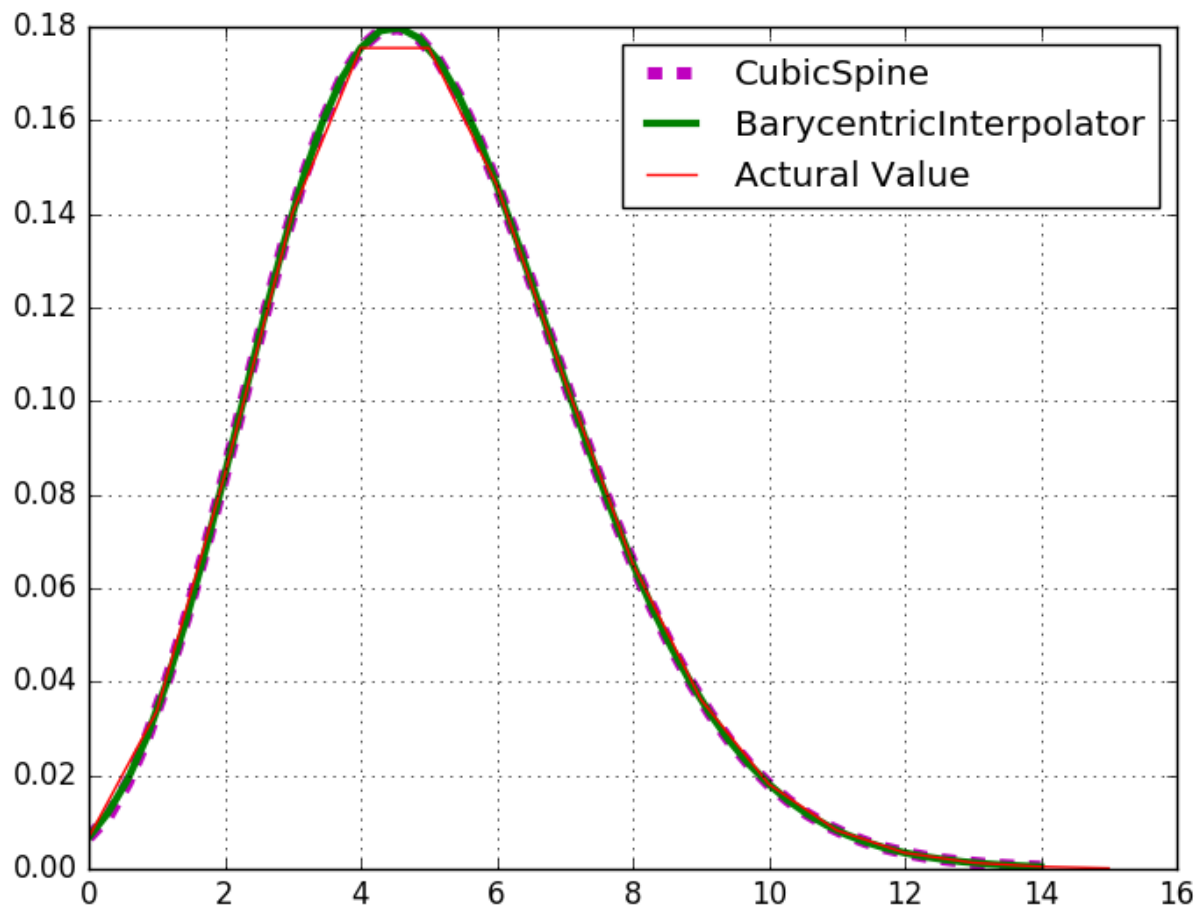
有：

$$f(x) = \frac{\sum_{j=0}^n \frac{u_j}{x - x_j} y_j}{\sum_{j=0}^n \frac{u_j}{x - x_j}}$$

□ 则函数  $f(x)$  在  $x_k$  处的值为  $y_k$ 。

■ 对于权值，可以选择  $\{u_j = (-1)^j, j = 0, 1, \dots, n\}$

# 样条插值 – 重心插值



# Demo

---

# 作业

---

- 实现任何一个函数曲线/曲面的Python显示。
  - Matplotlib
- 利用Python提供的SVD库函数，实现图像恢复。
- 数值计算

# 作业：数值计算

□ 对于某二分类问题，若构造了九个正确率都是0.6的分类器，采用少数服从多数的原则进行最终分类，则最终分类正确率是多少？

■ 若构造99个分类器呢？

```
def bagging(n, p):  
    p = 0.6  
    s = 0  
    for i in range(n / 2 + 1, n + 1):  
        s += c(n, i) * p ** i * (1 - p) ** (n - i)  
    return s  
  
if __name__ == "__main__":  
    for t in range(9, 100, 10):  
        print t, '次采样正确率: ', bagging(t, 0.6)
```

Ensemble

```
C:\Python27\python.exe D:/Python/Ensemble.py  
9 次采样正确率: 0.73343232  
19 次采样正确率: 0.813907978585  
29 次采样正确率: 0.863787051336  
39 次采样正确率: 0.897941368711  
49 次采样正确率: 0.922424437652  
59 次采样正确率: 0.940447995732  
69 次采样正确率: 0.953949756505  
79 次采样正确率: 0.964189692839  
89 次采样正确率: 0.972027516007  
99 次采样正确率: 0.97806955787
```

# 我们在这里

□ <http://wenda.ChinaHadoop.cn>

■ 视频/课程/社区

□ 微博

■ @ChinaHadoop

■ @邹博\_机器学习

□ 微信公众号

■ 小象学院

■ 大数据分析挖掘



---

感谢大家！

恳请大家批评指正！