

# 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



# 机器学习与数学分析

---



小象学院  
ChinaHadoop.cn

邹博

# 主要内容

---

- 机器学习与本课程示例概述
- 机器学习的角度看数学：
  - 数学分析
    - 导数与梯度
    - Taylor展式的应用
  - 概率论基础
    - 古典概型
    - 频率学派与贝叶斯学派
    - 常见概率分布
    - Sigmoid/Logistic函数的引入

# 什么是机器学习

- 对于某给定的任务T，在合理的性能度量方案P的前提下，某计算机程序可以自主学习任务T的经验E；随着提供合适、优质、大量的经验E，该程序对于任务T的性能逐步提高。
- 这里最重要的是机器学习的对象：
  - 任务Task, T，一个或者多个
  - 经验Experience, E
  - 性能Performance, P
- 即：随着任务的不断执行，经验的累积会带来计算机性能的提升。
  - Tom Michael Mitchell, 1997

# 换个表述

□ 机器学习是人工智能的一个分支。我们使用计算机设计一个系统，使它能够根据提供的训练数据按照一定的方式来学习；随着训练次数的增加，该系统可以在性能上不断学习和改进；通过参数优化的学习模型，能够用于预测相关问题的输出。

□ 思考：

■ 如何设计无人驾驶机动车？

# 无人驾驶汽车

- 汽车的无人汽车模块已经成熟：全自动公共交通工具已经出现在了世界上的多个城市。
  - Lutz探路者/CYCAB/Google
- 问题：如何设计自动驾驶系统？



# 人类的学习?



□ 如何从完全“无知”到掌握知识?

■ 语言/颜色/形状等特征统计

□ 有监督学习

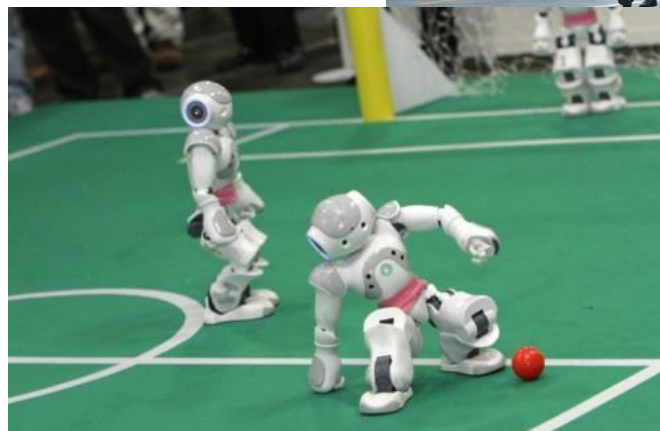
■ 月亮

□ 无监督学习

■ 9月3日：“阅兵”

□ 增强学习

■ 走路、踢球



# 机器学习的内涵与外延

## □ 机器学习可以解决什么

### ■ 给定数据的预测问题

□ 数据清洗/特征选择

□ 确定算法模型/参数优化

□ 结果预测

## □ 不能解决什么

### ■ 大数据存储/并行计算

### ■ 做一个机器人

## □ 举例：

### ■ 机器学习：“盯住2号位，她很容易起快球”

### ■ 传统算法：排球规则。





# 机器学习



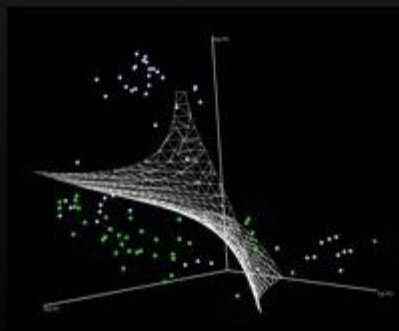
what society thinks I do

$$\begin{aligned}L_p &= \frac{1}{2} \|w\|^2 - \sum_{i,j} \alpha_{i,j} (x_i \cdot w + b) + \sum_i \alpha_i \\ \alpha_i &\geq 0, \forall i \\ w &= \sum_{i,j} \alpha_{i,j} x_i, \sum_i \alpha_i = 0 \\ \nabla \hat{g}(\theta_t) &= \frac{1}{n} \sum_{i=1}^n \nabla \ell(x_{i(t)}, y_{i(t)}; \theta_t) + \nabla r(\theta_t) \\ \theta_{t+1} &= \theta_t - \eta_t \nabla \ell(x_{i(t)}, y_{i(t)}; \theta_t) - \eta_t \cdot \nabla r(\theta_t) \\ \mathbb{E}_{i(t)}[\ell(x_{i(t)}, y_{i(t)}; \theta_t)] &= \frac{1}{n} \sum_i \ell(x_i, y_i; \theta_t).\end{aligned}$$

what other programmers think I do



what my friends think I do



what I think I do



what my parents think I do

```
>>> from sklearn import svm
```

what I really do

features

target

samples  
(train)

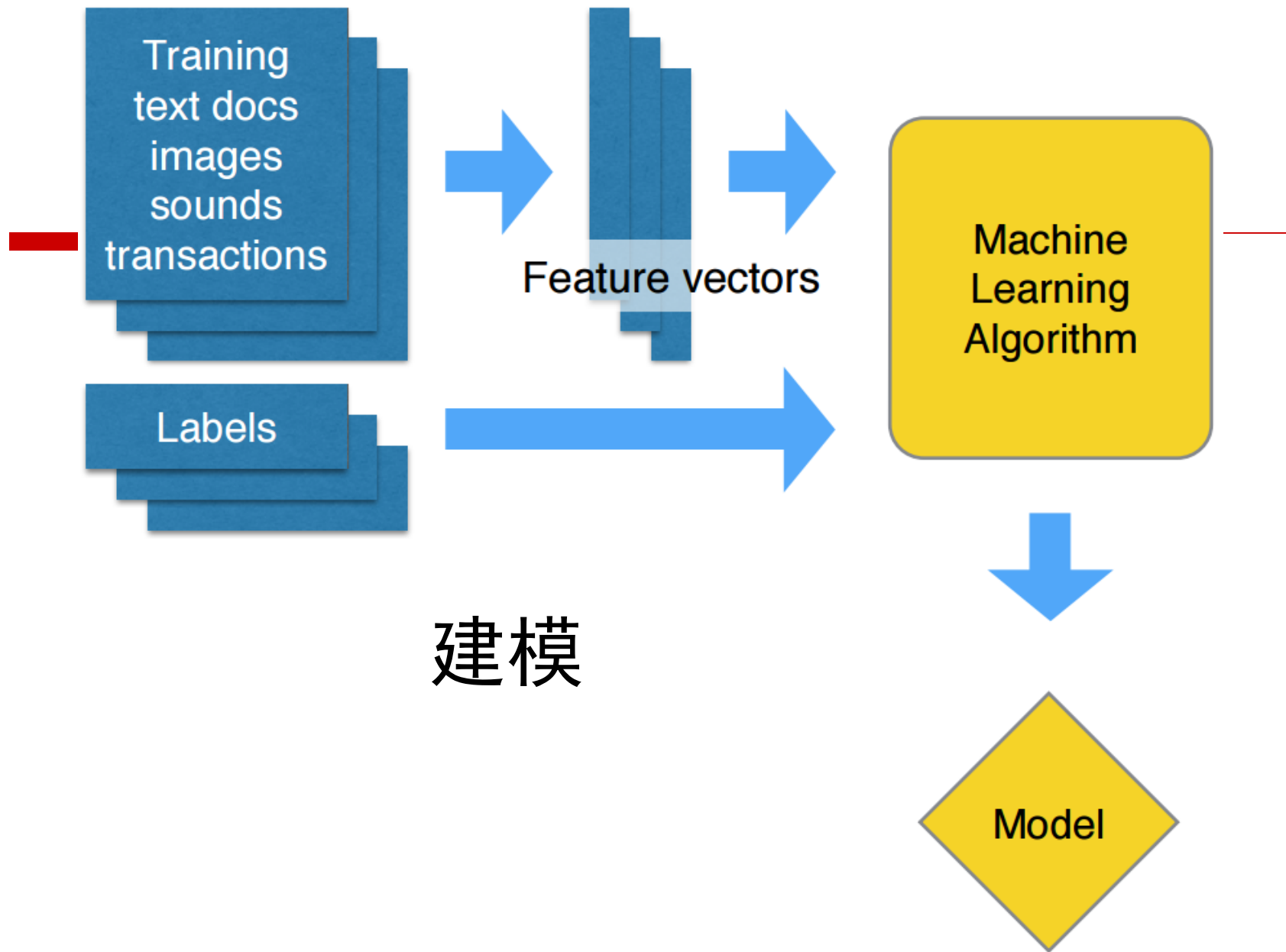
type (category)	# rooms (int)	surface (float m2)	public trans (boolean)
Apartment	3	50	TRUE
House	5	254	FALSE
Duplex	4	68	TRUE
Apartment	2	32	TRUE

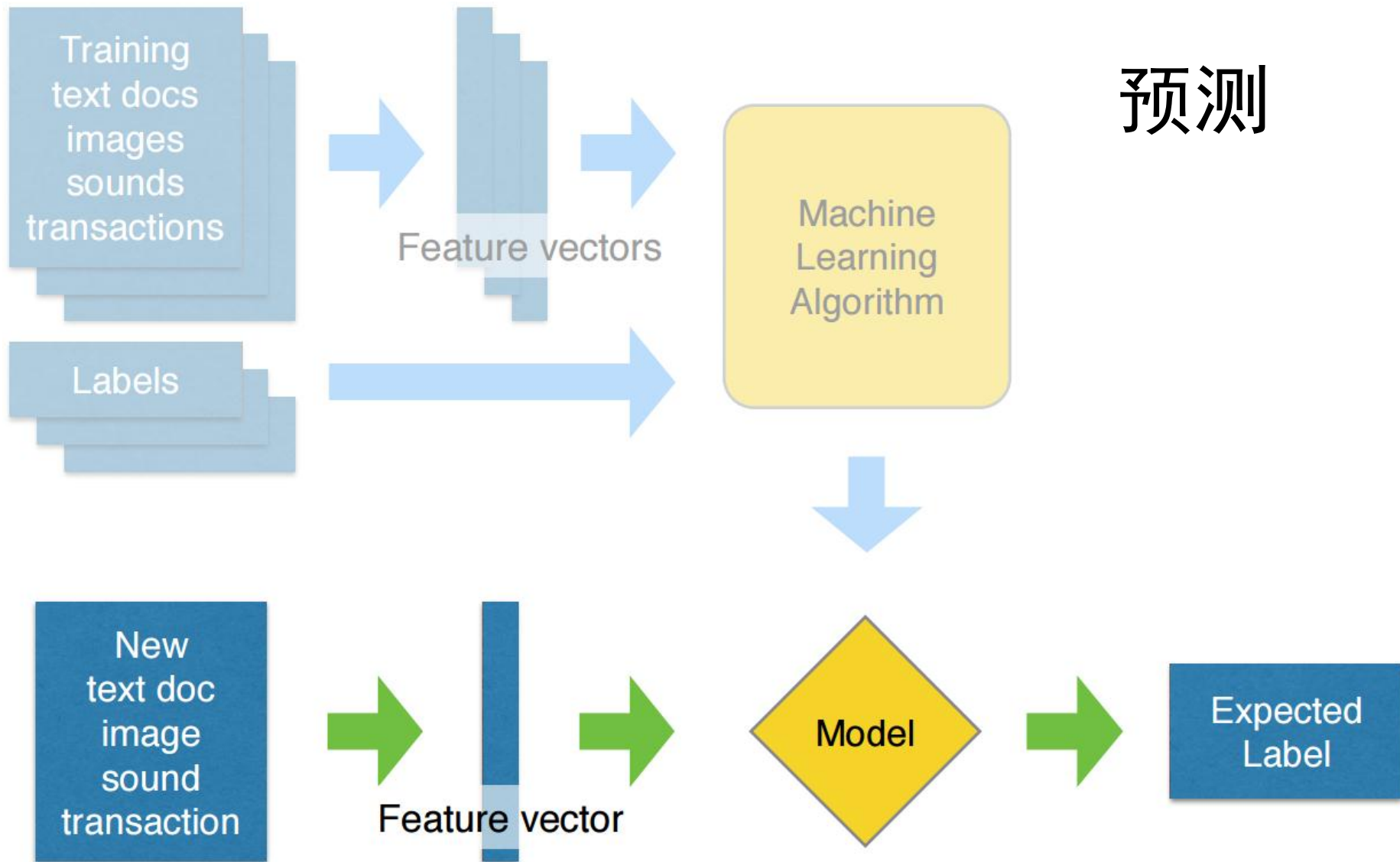
sold (float k€)
450
430
712
234

samples  
(test)

Apartment	2	33	TRUE
House	4	210	TRUE

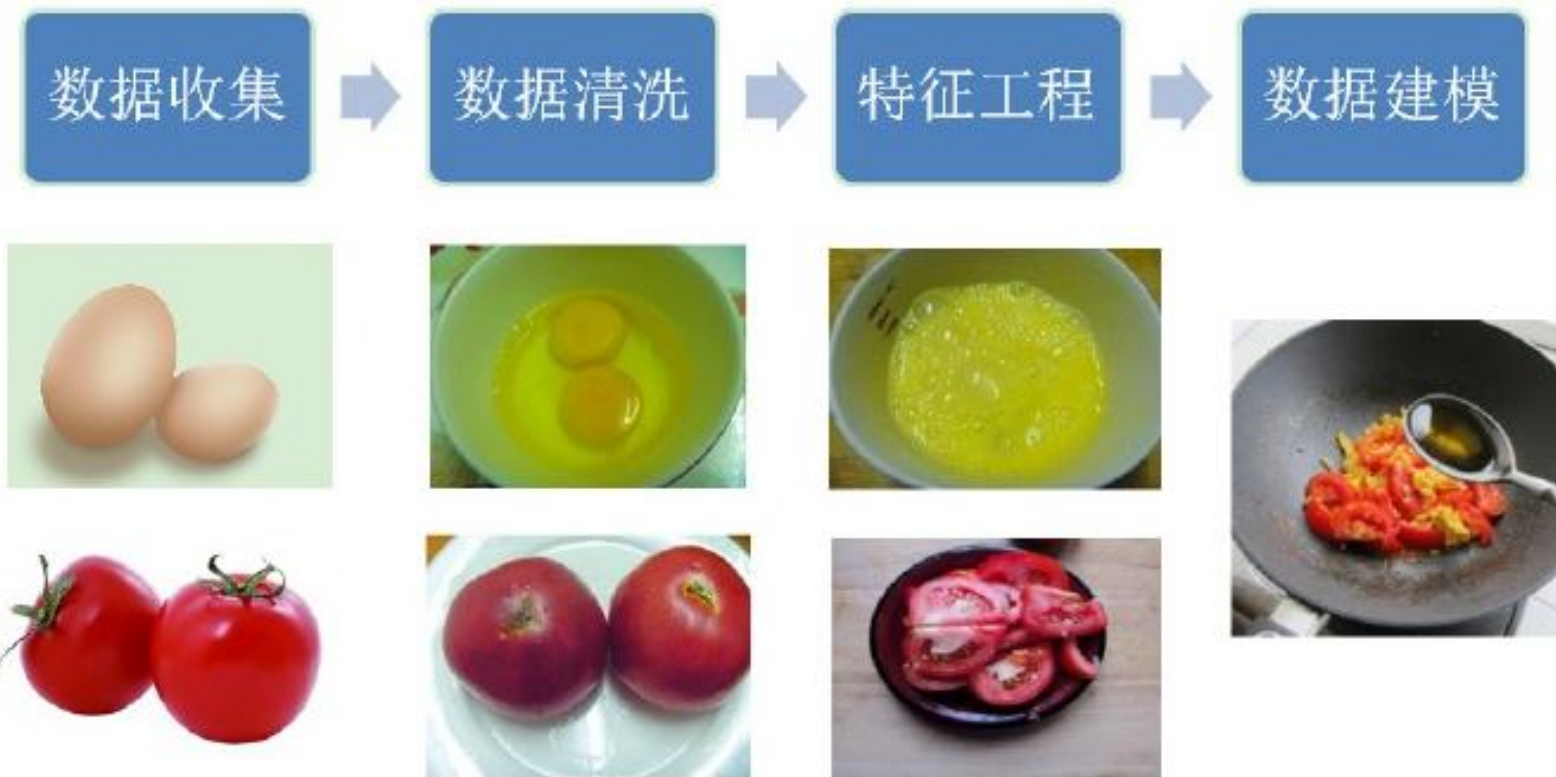
?
?



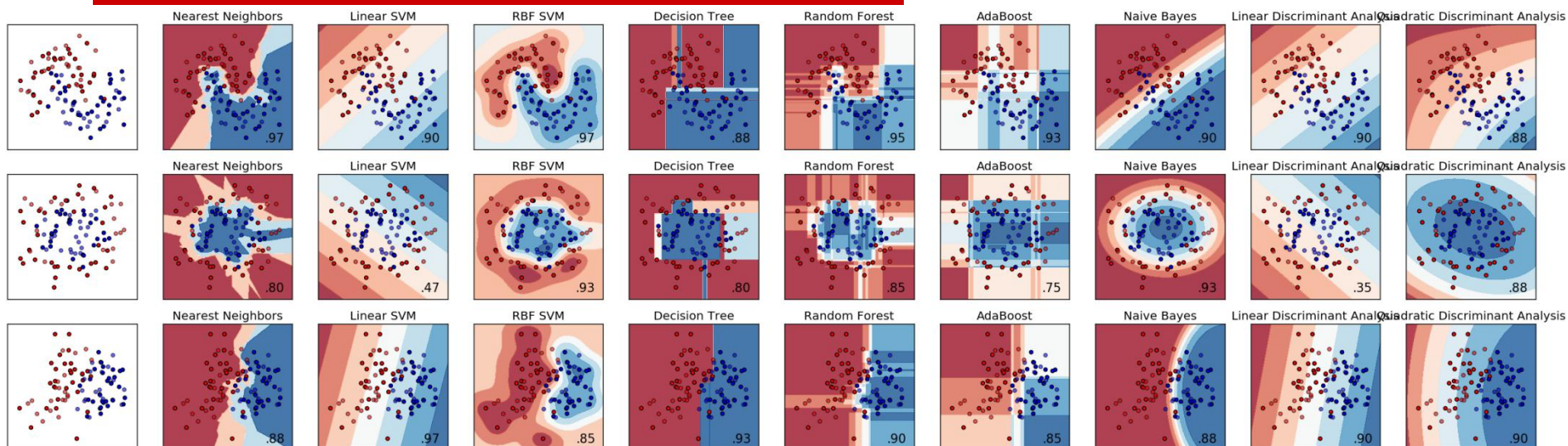


# 机器学习的一般流程

## 数据挖掘/机器学习的流程与西红柿炒鸡蛋



# 机器学习方法



- different assumptions on data
- different scalability profiles at training time
- different latencies at prediction time
- different model sizes (embedability in mobile devices)

# 思考：机器如何发现新词

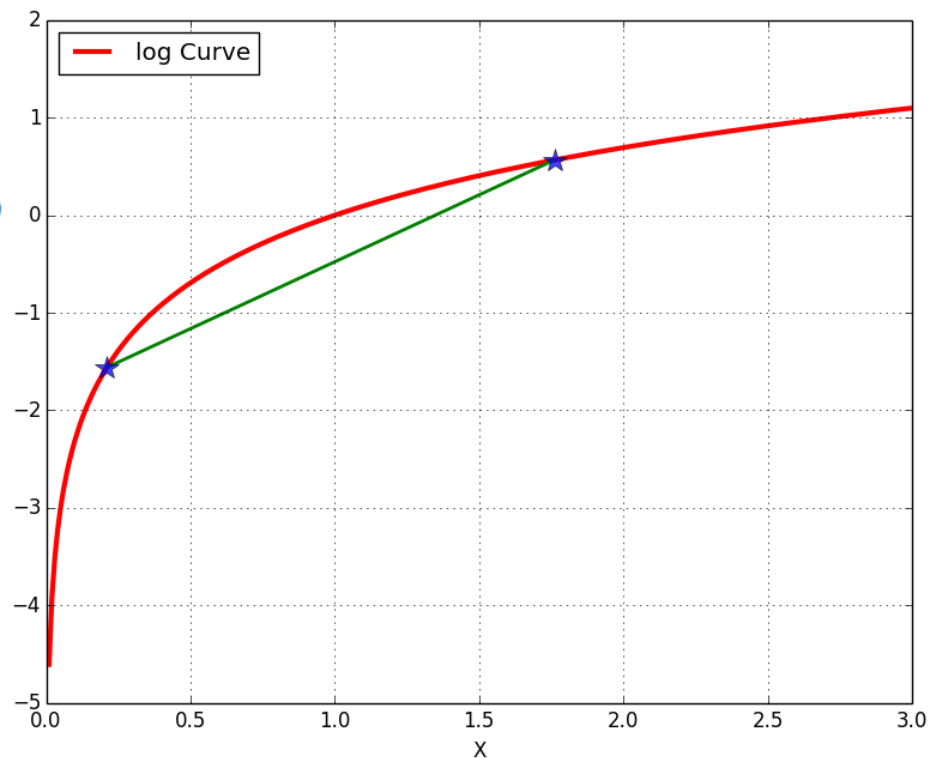
- 频数：Count(X)
- 凝固程度
  - $X = A.B$
  - $P(A)P(B)$  vs  $P(X)$
- 自由程度
  - $aXb$
  - 信息熵 $H(a)$ 、 $H(b)$
- 凝固程度和自由程度缺一不可。只考虑凝固程度，会找出“巧克”、“俄罗”、“颜六色”、“柴可夫”等“半个词”；只考虑自由程度，会把“吃了一顿”、“看了一遍”、“睡了一晚”、“去了一趟”中的“了一”提取出来，因为它的左右邻字都太丰富了。
  - 调参
- 问题：给定某长文本，如何利用上述参数设计可行算法？

# Python Code示例1

```
TestLog.py x
# -*- coding:utf8 -*-

import math
import matplotlib.pyplot as plt

if __name__ == "__main__":
    x = [float(i)/100.0 for i in range(1,300)]
    y = [math.log(i) for i in x]
    plt.plot(x, y, 'r-', linewidth=3, label='log Curve')
    a = [x[20], x[175]]
    b = [y[20], y[175]]
    plt.plot(a, b, 'g-', linewidth=2)
    plt.plot(a, b, 'b*', markersize=15, alpha=0.75)
    plt.legend(loc='upper left')
    plt.grid(True)
    plt.xlabel('X')
    plt.ylabel('log(X)')
    plt.show()
```

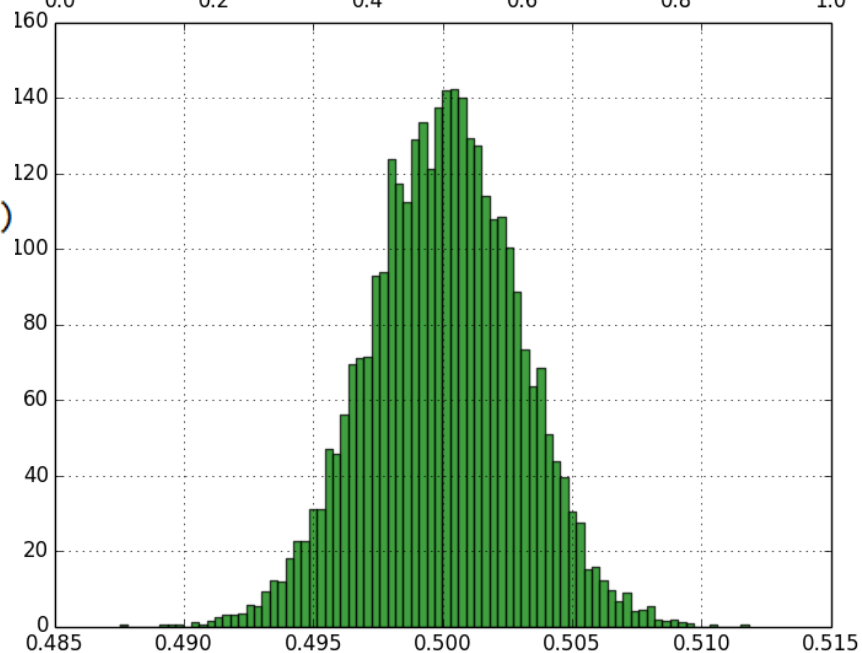
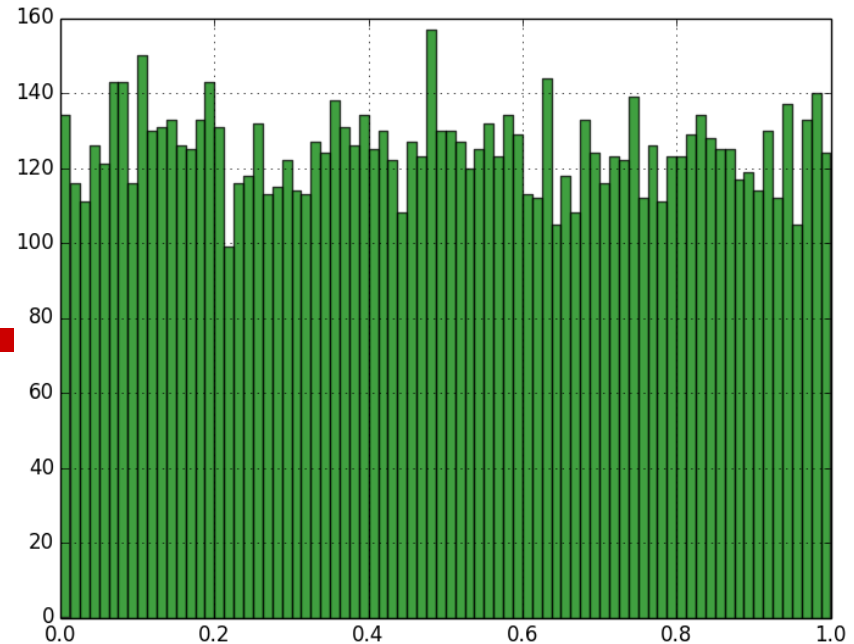




# Python Code示例2

```
if __name__ == "__main__":
    u = numpy.random.uniform(0.0, 1.0, 10000)
    plt.hist(u, 80, facecolor='g', alpha=0.75)
    plt.grid(True)
    plt.show()

    times = 10000
    for time in range(times):
        u += numpy.random.uniform(0.0, 1.0, 10000)
    print len(u)
    u /= times
    print len(u)
    plt.hist(u, 80, facecolor='g', alpha=0.75)
    plt.grid(True)
    plt.show()
```



# Py

脚本文件一般用 .py 后缀

coffeeghost-q-in-py.py x

中文用户一定得先用这行来声明编码,同时文件本身也得存储成UTF-8编码!

单行注释

导入其它代码模块

模块名,其实导入了 os.py

函数名"main"在这儿并不是必须的,调用在这段脚本的最后部分;

注意!Python最好也最个性的语法:  
使用缩进来代替其它语句块声明;  
一般建议每个层级用4个空格来缩进.

print 'Hello World!' 声明单行字符串,使用双/单引号都成,  
注意对字符串中的引号进行逃逸处理!

print "这是Alice\的问候."  
print '这是Bob\的问候.' 函数调用,声明在后述代码;

foo(5, 10) 字符可乘,等于:'====='

print '=' \* 10  
print '这将直接执行'+os.getcwd() 调用了os 模块中的函数

变量得先实例化  
才可进一步计算

counter = 0  
counter += 1 连接字符串  
内置的列表类型对象,其实可以包含不同类型数据,  
甚至可以包含其它列表对象;

单行的语句块,其实可以不换行的,  
但是,建议清晰起见,规范点:  
- 另起一行  
- 缩进一级

food = ['苹果', '杏子', '李子', '梨']  
for i in food:  
print '俺就爱整只:'+i 在循环中,i 指代了列表中按顺序的每个"food"

print '数到10'  
for i in range(10):  
print i range()内置函数,返回类似  
[0,1,2,3,4,5,6,7,8,9]  
的数字列表,注意 for in 循环语句使用冒号结束声明!

函数声明,  
注意使用冒号结束声明

def foo(param1, secondParam): 字符串的格式化输出基本类似C语言的

res = param1+secondParam  
print '%s 加 %s 等于 %s'%(param1, secondParam, res)

if res < 50:  
print '这个' 判定式也基本和C语言的相同

elif (res>=50) and ((param1==42) or (secondParam==24)):  
print '那个'

else:  
print '嗯...' 逻辑运算符,不使用 && 和 ||,使用直观的文单词

return res # 这是单行注释

'''这是多  
行注释.....''' 这都是合法注释

多行注释的内容不用遵守当前缩进  
只要开始的''' 缩进正确就成!

每级语法块不用}之类的括号引领!  
直接回车+4空格  
(当然,要在当前缩进基础上)

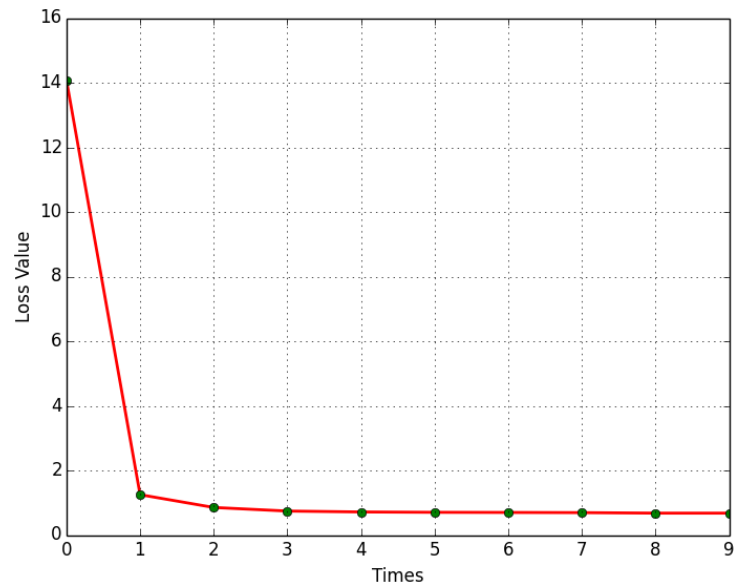
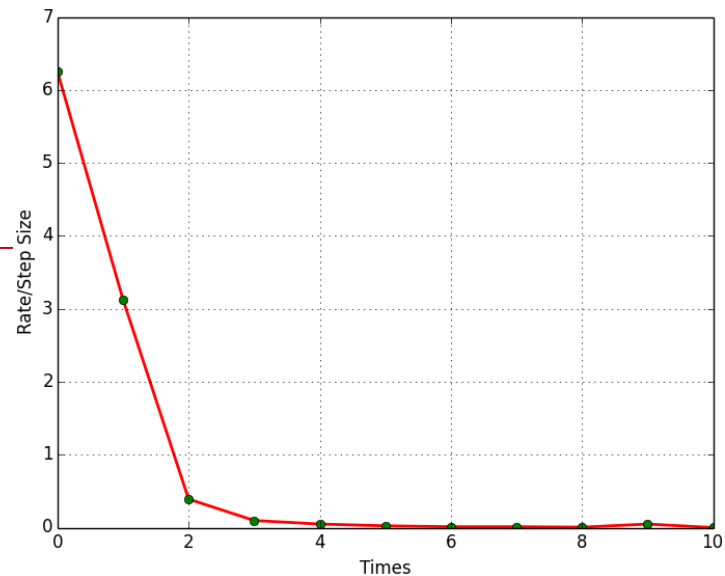
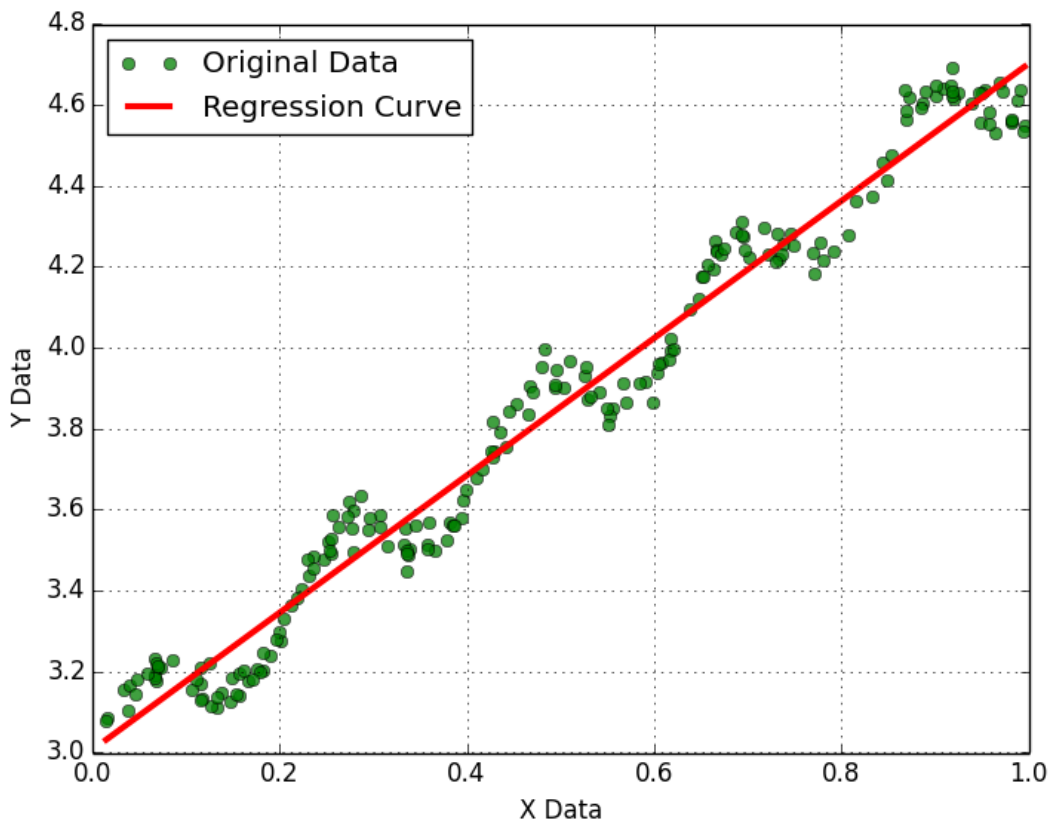
if \_\_name\_\_ == '\_\_main\_\_':  
main() 一般在脚本最后调用主函数 main();而且使用 内置的运行脚本名来判定;  
当且仅当我们直接运行当前脚本时, \_\_name\_\_ 才为 \_\_main\_\_  
这样当脚本被当作模块进行 import 导入时,并不运行 main()  
所以,一般这里是进行测试代码安置的...

在Ms中很好的支持UTF-8的编辑器不多,  
跨平台又支持Py特性的更少;  
好在我们有 Limodou 贡献的 UliPad  
这一编辑器本身就是Py实现的!  
(基于wxPython)

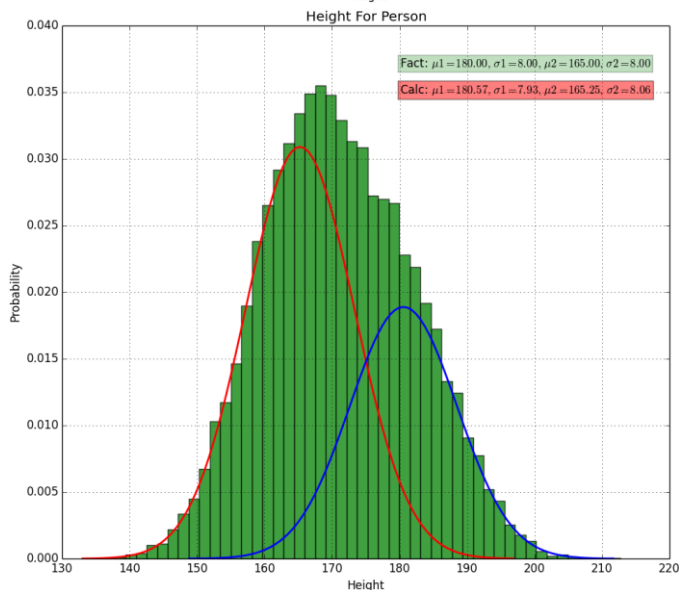
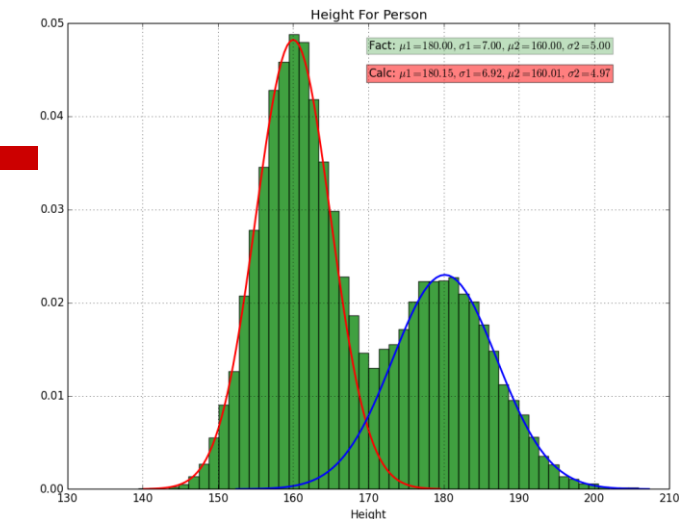


用冒号来结束判断句,  
在 if elif else 行最后

# 线性回归、rate、Loss



# EM Code



```
em3.py x
def calcEM(height):
    N = len(height)
    gp = 0.5 #girl probability
    bp = 0.5 #boy probability
    gm,gsigma = min(height),1 #先验: 直接取最大和最小值
    bm,bsigma = max(height),1
    ggamma = range(N)
    bgamma = range(N)
    cur = [gp, bp, gm, gsigma, bm, bsigma]
    now = []

    times = 0
    while times < 100:
        i = 0
        for x in height:
            ggamma[i] = gp * gauss(x, gm, gsigma)
            bgamma[i] = bp * gauss(x, bm, bsigma)
            s = ggamma[i] + bgamma[i]
            ggamma[i] /= s
            bgamma[i] /= s
            i += 1

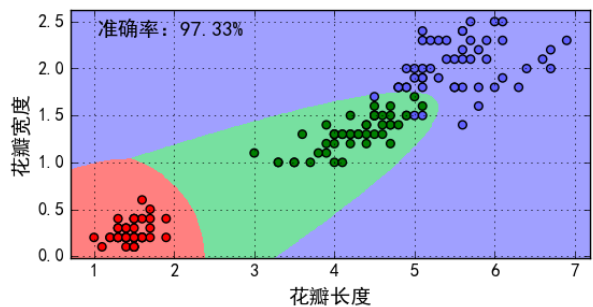
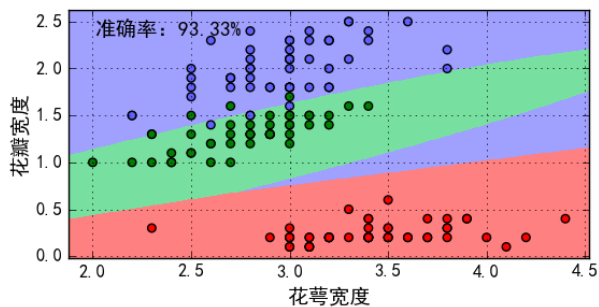
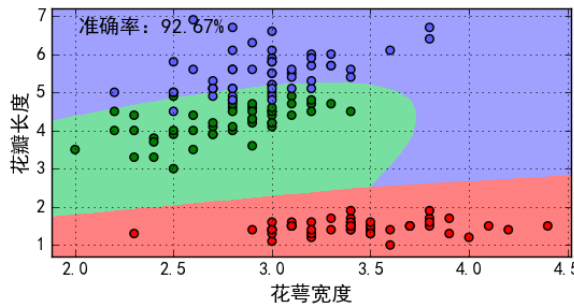
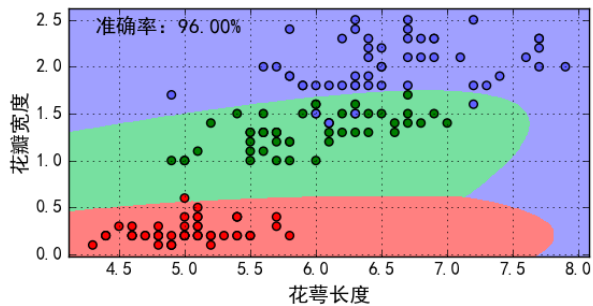
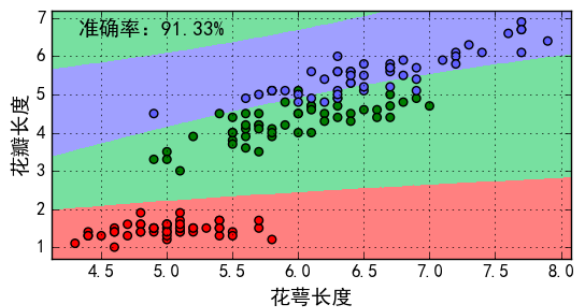
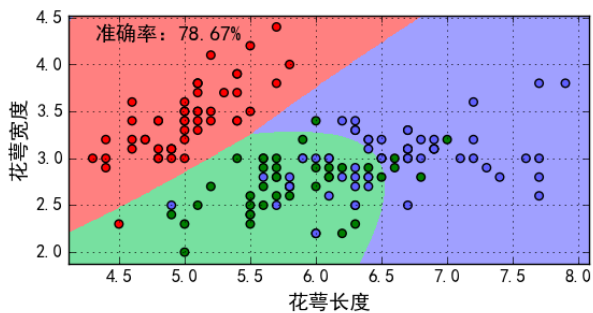
        gn = sum(ggamma)
        gp = float(gn) / float(N)
        bn = sum(bgamma)
        bp = float(bn) / float(N)
        gm = averageweight(height, ggamma, gn)
        gsigma = varianceWeight(height, ggamma, gm, gn)
        bm = averageweight(height, bgamma, bn)
        bsigma = varianceWeight(height, bgamma, bm, bn)

        now = [gp, bp, gm, gsigma, bm, bsigma]
        if isSame(cur, now):
            break
        cur = now
        print "Times:\t", times
        print "Girl mean/gsigma:\t", gm,gsigma
        print "Boy mean/bsigma:\t", bm,bsigma
        print "Boy/Girl:\t", bn, gn, bn+gn
        print "\n\n"
        times += 1

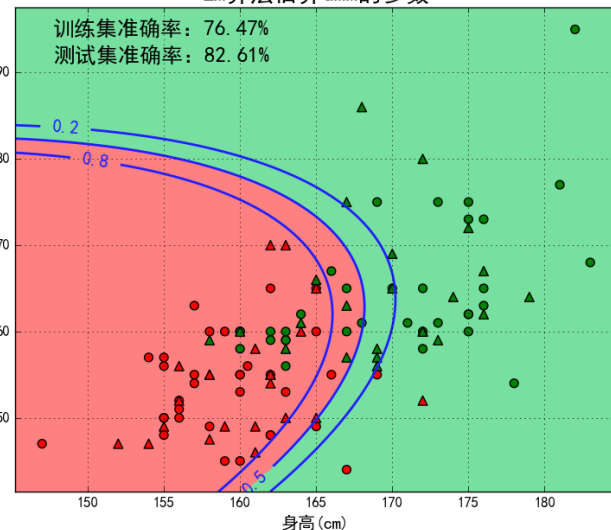
    return now
```

# EM算法

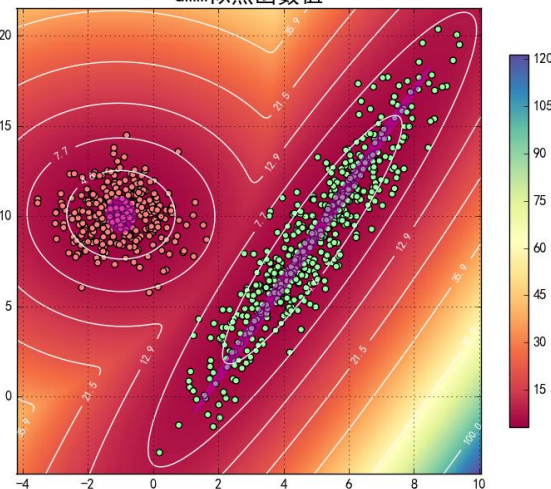
EM算法无监督分类鸢尾花数据



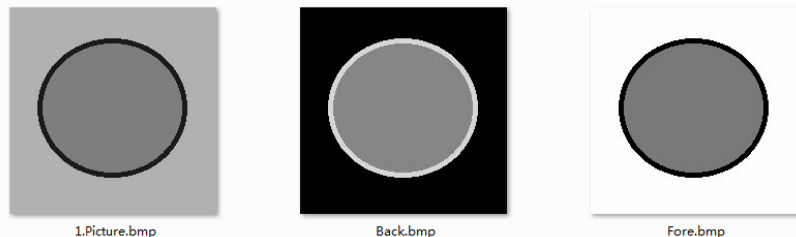
EM算法估算GMM的参数



GMM似然函数值



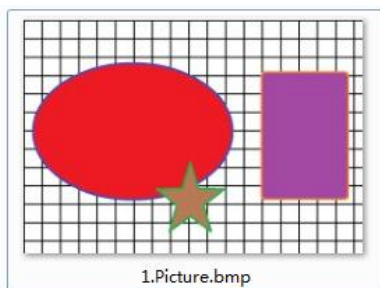
# GMM与图像



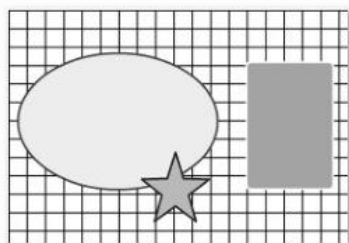
1.Picture.bmp

Back.bmp

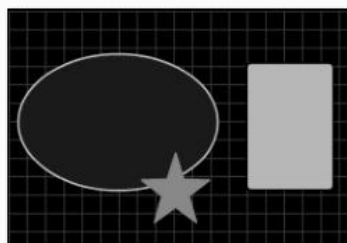
Fore.bmp



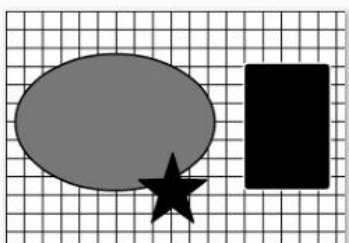
1.Picture.bmp



2.Red Channel.bmp



back.bmp



fore.bmp

```
def composite(band, parameter):  
    c1 = parameter[0]  
    mu1 = parameter[2]  
    sigma1 = parameter[3]  
    c2 = parameter[1]  
    mu2 = parameter[4]  
    sigma2 = parameter[5]  
  
    p1 = []  
    p2 = []  
    for pixel in band:  
        p1.append(c1 * gauss(pixel, mu1, sigma1))  
        p2.append(c2 * gauss(pixel, mu2, sigma2))  
  
    scale(p1) #灰度均衡  
    scale(p2)  
    return [p1, p2]  
  
if __name__ == "__main__":  
    im = Image.open('.\\Pic\\test.bmp')  
    print im.format, im.size, im.mode  
  
    im = im.split()[0] #只处理第一个通道  
    nb = [] #处理后的新通道  
    data = list(im.getdata())  
    parameter = GMM(data)  
    t = composite(data, parameter)  
  
    im1 = Image.new('L', im.size)  
    im1.putdata(t[0])
```

# 图像的卷积



laplacian.png



laplacian2.png



prewitt.png



prewitt\_x.png



prewitt\_y.png



soble.png



soble\_x.png



soble\_y.png

# 去均值ICA分离

源信号1

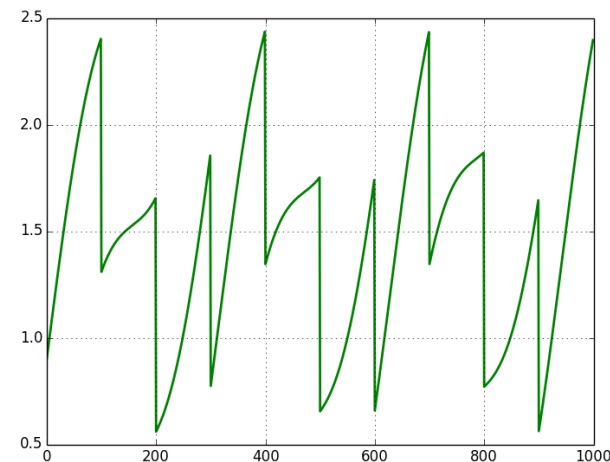
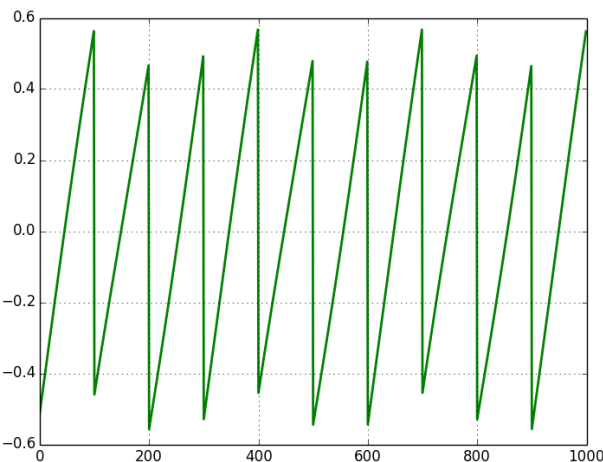
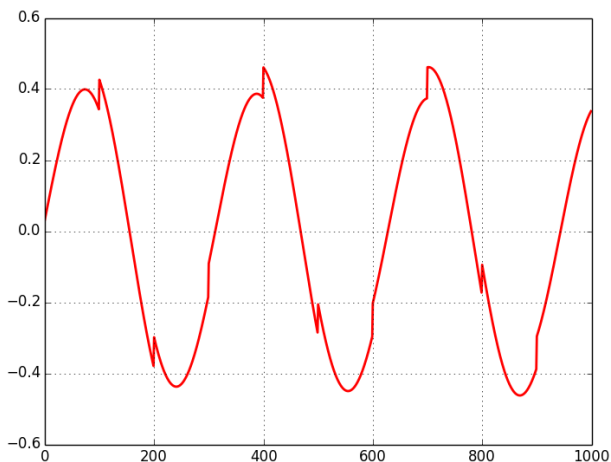
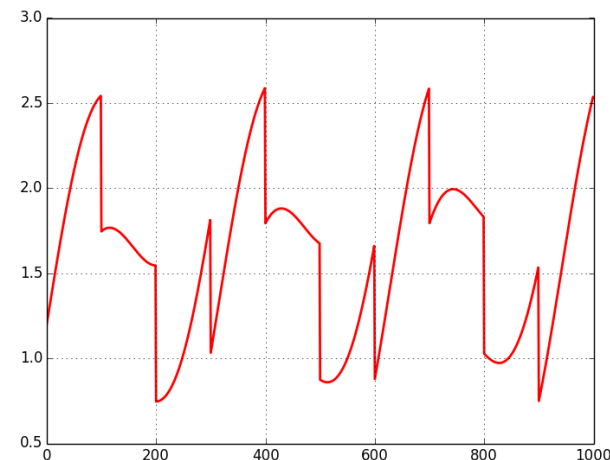
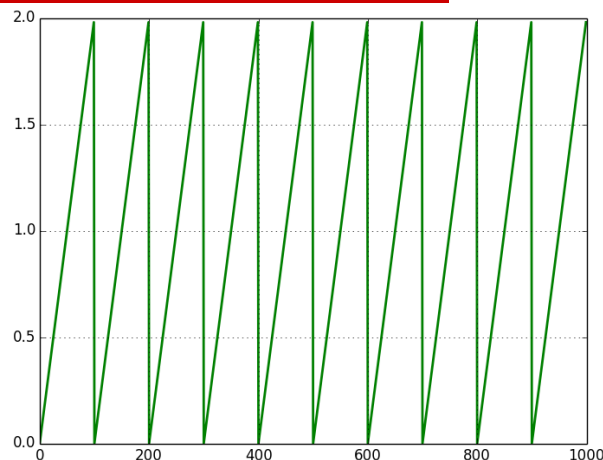
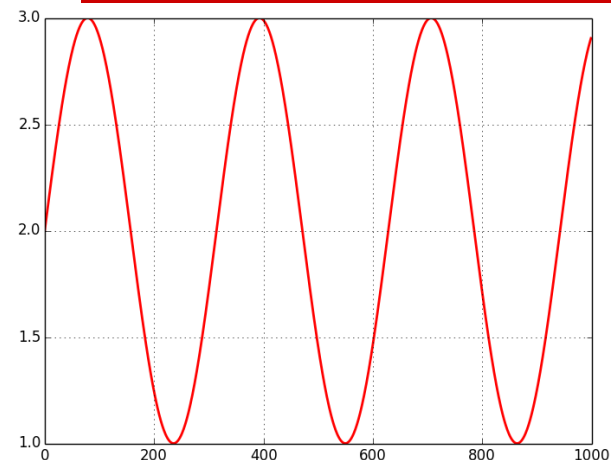
源信号2

混合信号1

独立成分1

独立成分2

混合信号2





# 带噪声的信号分离

源信号1

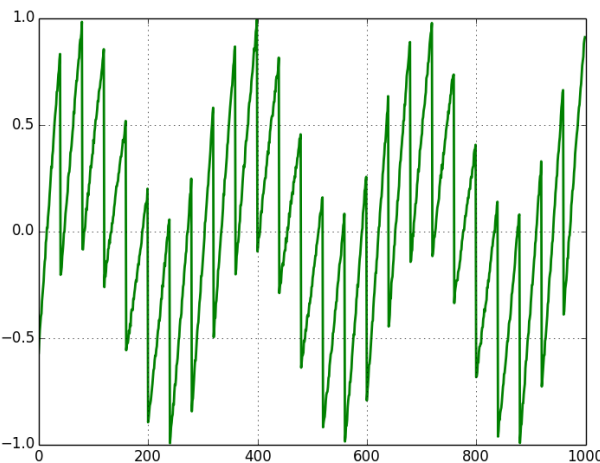
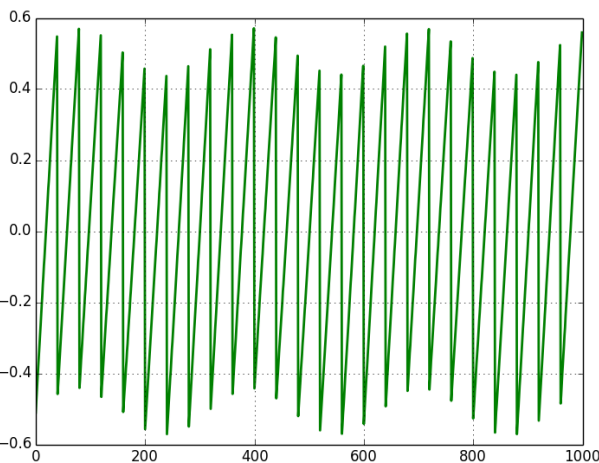
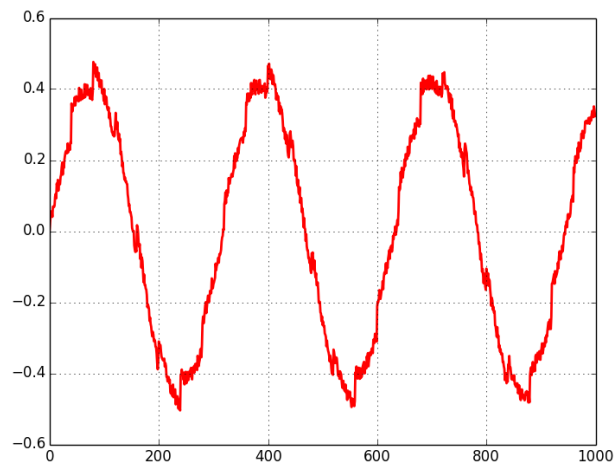
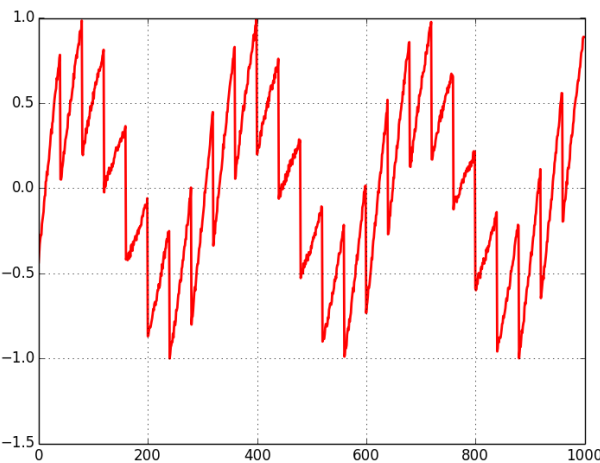
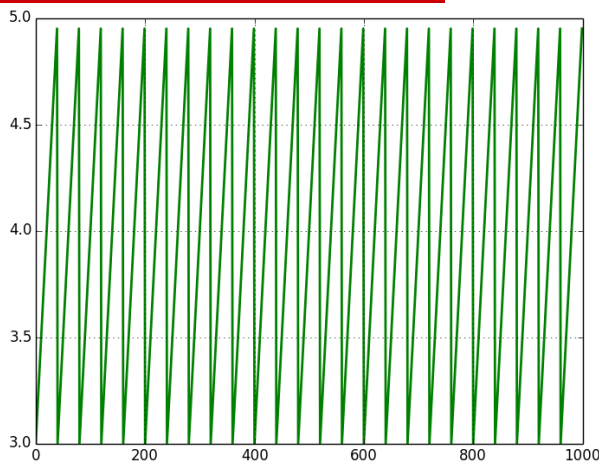
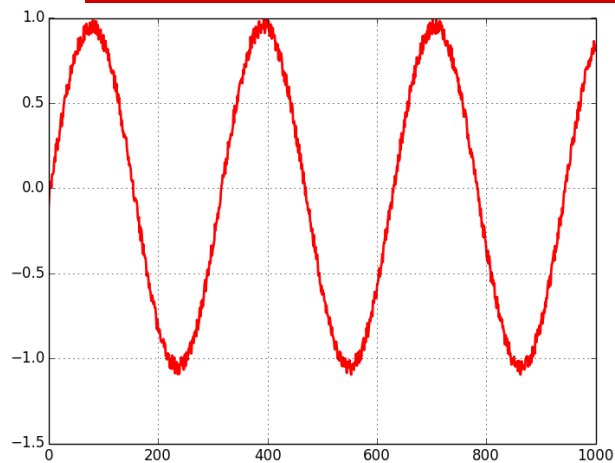
源信号2

混合信号1

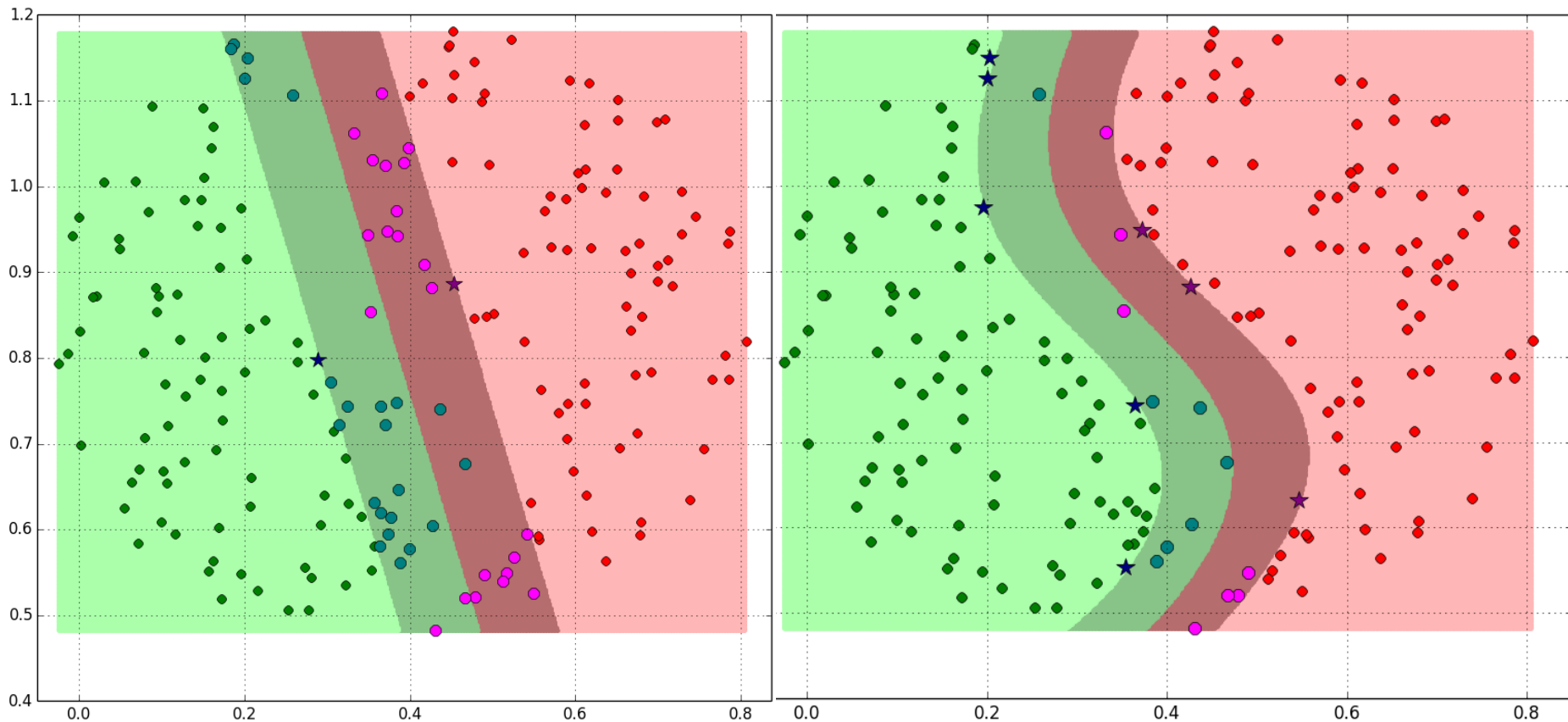
独立成分1

独立成分2

混合信号2



# SVM: 高斯核函数的影响

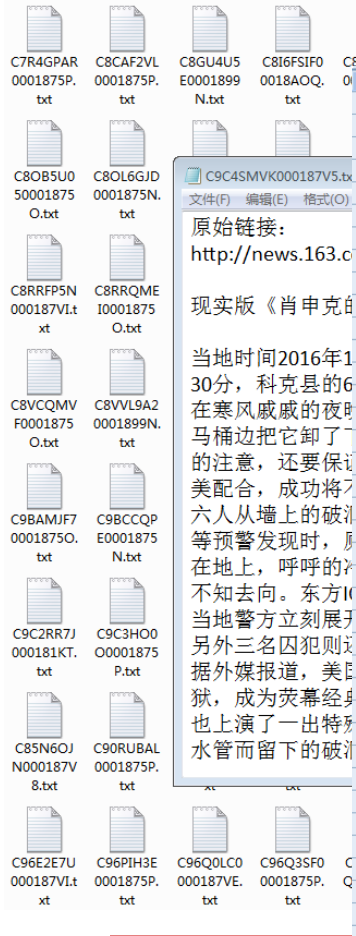


# Crawler爬取数据

```

writer.writerow(['区域', '小区名称', '户型', '面积', '价格(万)', '单价(元/平米)',
                '性质', '朝向', '装修', '是否有电梯', '楼层', '建筑年代', '楼型'])
res = requests.get('http://bj.lianjia.com/ershoufang/')
res = res.text.encode(res.encoding).decode('utf-8')
soup = BeautifulSoup(res, 'html.parser')
# print soup.prettify()
districts = soup.find(name='div', attrs={'data-role': 'ershoufang'}) # <div data-role="ers
for district in districts.find_all(name='a'):
    print district['title']
    district_name = district.text # '东城', '西城', '朝阳', '海淀'.....
    url = '%s%s' % (url_main, district['href'])
    print url
    res = requests.get(url)
    res = res.text.encode(res.encoding).decode('utf-8')

```



	A	B	C	D	E	F	G	H	I	J	K	L	M
1	区域	小区名称	户型	面积	价格(万)	单价(元/平米)	性质	朝向	装修	是否有电梯	楼层	建筑年代	楼型
2	东城	桃杨路北里12号院	3室1厅	98.17	638	64990	房本满五年	东西	简装	无电梯	中楼层(共6层)	1999	板楼
3	东城	民安小区东扬威街	3室1厅	103.29	1000	96815	房本满五年	西南北	简装	有电梯	低楼层(共16层)	2002	塔楼
4	东城	民安小区东直门内北小街	2室1厅	74.54	732	98203	房本满五年	南北	简装	有电梯	高楼层(共16层)	2003	板楼
5	东城	海运仓小区	2室1厅	64.82	700	107992	房本满五年	南北	精装	无电梯	中楼层(共7层)	2003	板楼
6	东城	新景家园西区	1室1厅	57.18	560	97937	房本满五年	西南	精装	有电梯	中楼层(共17层)	2004	板塔结合
7	东城	小羊宜宾胡同	4室1厅	93.81	910	97005	房本满五年	南北	简装	有电梯	低楼层(共16层)	1992	塔楼
8	东城	禾风相府	3室2厅	163.83	1801	109932	房本满五年	南北	其他	有电梯	低楼层(共14层)	2008	板塔结合
9	东城	西革新里110号院	2室1厅	61.4	400	65147	房本满五年	南北	简装	无电梯	高楼层(共9层)	1992	板楼
10	东城	东中街	4室1厅	200.77	1770	88161	房本满五年	南北	精装	有电梯	中楼层(共6层)	1999	板塔结合
11	东城	东花市北里中区	3室2厅	98.65	990	100355	房本满五年	南北	精装	无电梯	中楼层(共5层)	1994	板楼
12	东城	民安小区东直门北中街	2室1厅	75.11	722	96126	房本满五年	东西	精装	有电梯	中楼层(共11层)	2002	板塔结合
13	东城	中景濠庭	1室1厅	74.21	630	84895	房本满五年	北	精装	有电梯	低楼层(共22层)	2001	板塔结合
14	东城	华府景园	3室2厅	149.97	1350	90019	房本满五年	东西	简装	有电梯	中楼层(共12层)	2001	板塔结合
15	东城	朝阳门南小街	2室1厅	69.39	650	93674	房本满五年	东西	简装	有电梯	低楼层(共8层)	2004	板塔结合
16	东城	兴隆都市馨园	3室2厅	130.7	980	74981	房本满五年	南北	精装	无电梯	高楼层(共6层)	2003	板塔结合
17	东城	保利蔷薇	2室2厅	88.63	960	108316	房本满五年	西北	精装	有电梯	中楼层(共16层)	2008	板塔结合
18	东城	小羊宜宾胡同	2室1厅	61.77	650	105230	房本满五年	南	简装	有电梯	低楼层(共16层)	1992	塔楼
19	东城	安化北里	3室1厅	80.34	730	90864	房本满五年	西南	简装	有电梯	中楼层(共17层)	1989	塔楼
20	东城	香饵胡同	2室1厅	67.92	815	119995	房本满五年	南北	精装	无电梯	低楼层(共6层)	2002	板楼
21	东城	什锦花园胡同21号院	2室1厅	48	650	135417	房本满五年	南北	精装	无电梯	高楼层(共5层)	1985	板楼
22	东城	富贵园一区	3室2厅	144.45	1650	114227	房本满五年	南北	精装	有电梯	低楼层(共11层)	2003	板楼
23	东城	富莱茵花园	3室1厅	98.38	570	57939	房本满五年	西北	简装	有电梯	中楼层(共21层)	1995	塔楼
24	东城	丽水湾畔家园	2室1厅	130	1100	84616	房本满五年	西	简装	有电梯	高楼层(共21层)	2002	塔楼
25	东城	国瑞城中区	3室1厅	90.47	880	97270	房本满五年	西南	精装	有电梯	低楼层(共16层)	2006	板塔结合
26	东城	金鱼池中区	4室1厅	139.34	850	61002	房本满五年	南北	精装	无电梯	低楼层(共6层)	2002	板楼
27	东城	新景家园东区	3室1厅	92.82	1030	110968	房本满五年	南	精装	有电梯	高楼层(共13层)	2002	板塔结合
28	东城	本家润园三期	3室1厅	97.45	930	95434	房本满五年	南	精装	有电梯	高楼层(共18层)	2005	板塔结合
29	东城	新景家园东区	1室1厅	59.68	595	99699	房本满五年	南	简装	有电梯	中楼层(共17层)	2002	板楼
30	东城	北京上舍	1室0厅	61.2	580	94772	房本满五年	北	简装	有电梯	高楼层(共15层)	2008	塔楼
31	东城	和平里七区	2室1厅	41.82	450	107605	房本满五年	南	简装	无电梯	高楼层(共5层)	1983	板楼
32	东城	东土城路13号院	1室0厅	51.15	372	72728	房本满两年	南	精装	无电梯	中楼层(共5层)	1994	板楼
33	东城	金世纪嘉园	2室1厅	116.89	950	81273	房本满五年	东南	精装	有电梯	高楼层(共21层)	2003	塔楼
34	东城	民安小区民安街	2室1厅	74.29	868	116840	房本满五年	南北	精装	无	低楼层(共6层)	2002	板楼
35	东城	金鱼池中区	4室2厅	135.91	950	69900	房本满五年	南北	简装	无电梯	高楼层(共6层)	2002	板塔结合
36	东城	海晟名苑	2室1厅	140.53	1550	110297	房本满五年	东西	简装	有电梯	低楼层(共20层)	2003	板塔结合
37	东城	京城仁合	2室2厅	114.51	1010	88202	房本满五年	西南北	精装	有电梯	高楼层(共16层)	2002	板塔结合
38	东城	西水井胡同	2室1厅	74.95	680	90728	房本满五年	东西	精装	有电梯	中楼层(共16层)	2003	板塔结合
39	东城	新奥洋房	2室1厅	93.76	698	74446	房本满五年	东	精装	有电梯	中楼层(共11层)	2005	板塔结合
40	东城	西水井胡同	2室1厅	74.95	690	92062	房本满五年	东西	简装	有电梯	中楼层(共16层)	2003	板塔结合
41	东城	仓南胡同	2室1厅	77.8	920	118252	房本满五年	东西	简装	无电梯	中楼层(共6层)	2003	板楼
42	东城	新奥洋房	2室1厅	123.21	910	73858	房本满五年	东南西	精装	有电梯	高楼层(共13层)	2006	板塔结合
43	东城	北河沿大街	1室1厅	41.47	490	118158	房本满两年	南	简装	无电梯	高楼层(共6层)	1987	板楼
44	东城	东直门内大街	2室2厅	80.52	880	109290	房本满五年	东南	精装	有电梯	高楼层(共16层)	2003	塔楼

# HMM分词: MLE

## 前言

数据，数据，数据！想必在新闻、报刊等媒介的持续冲击下，人们无法摆脱大的洗礼。现实需求推动了对数据的学习。这些数据来自于社交媒体、智能手机、“物联网”、传感器等任何可以产生数据的多数据挖掘的宣传着重于数据规模数据洪水(data flood)的预言告诉人们我们无法实时处理这些数据，硬件推销人员会进一步卖给我们需要的服务，以期能够满足处理速度的要求。从某种程度上来说，他们是错的，但是我们值得停下来思考片刻，并对手边的任务进行适当的再认识。

近年来，数据挖掘和机器学习在我们周围持续火爆，各种媒体也不断推送着海量的数据。仔细观察就能发现，实际应用中的那些机器学习算法与多年前并没有什么两样；它们只是在应用的数据规模上有些不同。历数一下产生数据的组织，至少在我看来，数目其实并不多。无非是Google、Facebook、Twitter、Netflix以及其他为数不多的机构在使用若干学习算法和工具，这些算法和工具使得他们能够对数据进行测试分析。那么，真正的问题是：“对于其他人，大数据框架下的算法和工具的作用是什么呢？”

我承认本书将多次提及大数据和机器学习之间的关系，这是我无法忽视的一个客观问题；但是它只是一个很小的因素，终极目标是如何利用可用数据获取数据的本质

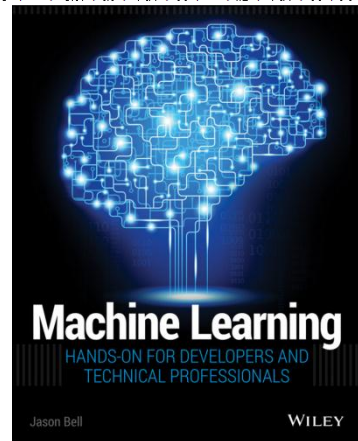
```
if __name__ == "__main__":
    pi, A, B = load_train()
    f = file("../text\\novel.txt")
    data = f.read()[3:].decode('utf-8')
    f.close()
    decode = viterbi(pi, A, B, data)
    segment(data, decode)
```

bug HMM

Console Frames Variables Watches

Connected to pydev debugger (build 139.1001)

我 | 与 | 地 | 坛 |  
史 | 铁 | 生 |  
— |  
我 | 在 | 好 | 几 | 篇 | 小 | 说 | 中 | 都 | 提 | 到 | 过 | 一 | 座 | 废 | 弃 | 的 | 古 | 园 |， | 实 | 际 | 就 | 是 | 地 | 坛 |。 | 许 | 多 | 年 | 前 | 旅 | 游 | 业 | 还 | 没 | 有 | 开 | 地 | 坛 | 离 | 我 | 家 | 很 | 近 |。 | 或 | 者 | 说 | 我 | 家 | 离 | 地 | 坛 | 很 | 近 |。 | 总 | 之 |， | 只 | 好 | 认 | 为 | 这 | 是 | 缘 | 分 |。 | 地 | 坛 | 在 | 我 | 出 | 生 | 前 | 它 | 等 | 待 | 我 | 出 | 生 |， | 然 | 后 | 又 | 等 | 待 | 我 | 活 | 到 | 最 | 狂 | 妄 | 的 | 年 | 龄 | 上 | 忽 | 地 | 残 | 废 | 了 | 双 | 腿 |。 | 四 | 百 | 多 | 年 | 里 |， | 它 | 一 | 面 | 自 | 从 | 那 | 个 | 下 | 午 | 我 | 无 | 意 | 中 | 进 | 了 | 这 | 园 | 子 |， | 就 | 再 | 没 | 长 | 久 | 地 | 离 | 开 | 过 | 它 |。 | 我 | 一 | 下 | 子 | 就 | 理 | 解 | 了 | 它 | 的 | 两 | 条 | 腿 | 残 | 废 | 后 | 的 | 最 | 初 | 几 | 年 |， | 我 | 找 | 不 | 到 | 工 | 作 |， | 找 | 不 | 到 | 去 | 路 |， | 忽 | 然 | 间 | 几 | 乎 | 什 | 么 | 都 | 找 | 不 | 到 | 了 |， | 除 | 去 | 几 | 座 | 殿 | 堂 | 我 | 无 | 法 | 进 | 去 |， | 除 | 去 | 那 | 座 | 祭 | 坛 | 我 | 不 | 能 | 上 | 去 | 而 | 只 | 能 | 从 | 各 | 个 | 角 | 度 | 张 | 望 | 它 |， | 地 | 坛 | 的 | 刺 | 下 | 的 | 就 | 是 | 怎 | 样 | 活 | 的 | 问 | 题 | 了 |， | 这 | 却 | 不 | 是 | 在 | 某 | 一 | 个 | 瞬 | 间 | 就 | 能 | 完 | 全 | 想 | 透 | 的 |、 | 不 | 是 | 一 | 次 | 性 | 能 | 二 |  
我 | 才 | 想 | 到 |， | 当 | 年 | 我 | 总 | 是 | 独 | 自 | 跑 | 到 | 地 | 坛 | 去 |， | 曾 | 经 | 给 | 母 | 亲 | 出 | 了 | 一 | 个 | 怎 | 样 | 的 | 问 | 题 |。 |  
她 | 不 | 是 | 那 | 种 | 光 | 会 | 疼 | 爱 | 儿 | 子 | 而 | 不 | 懂 | 得 | 理 | 解 | 儿 | 子 | 的 | 母 | 亲 |。 | 她 | 知 | 道 | 我 | 心 | 里 | 的 | 苦 | 闷 |， | 知 | 道 | 不 | 该 | 阻 | 止 | 有 | 一 | 回 | 我 | 摇 | 车 | 出 | 了 | 小 | 孩 |； | 想 | 起 | 一 | 件 | 什 | 么 | 事 | 又 | 返 | 身 | 回 | 来 |， | 看 | 见 | 母 | 亲 | 仍 | 站 | 在 | 原 | 地 |， | 还 | 是 | 送 | 我 | 有 | 一 | 次 | 与 | 一 | 个 | 作 | 家 | 朋 | 友 | 聊 | 天 |， | 我 | 问 | 他 | 学 | 写 | 作 | 的 | 最 | 初 | 动 | 机 | 是 | 什 | 么 |？ | 他 | 想 | 了 | 一 | 会 | 说 |： | “ 为 | 我 | 母 | 亲 |。 | 在 | 我 | 的 | 头 | 一 | 篇 | 小 | 说 | 发 | 表 | 的 | 时 | 候 |， | 在 | 我 | 的 | 小 | 说 | 第 | 一 | 次 | 获 | 奖 | 的 | 那 | 些 | 日 | 子 | 里 |， | 我 | 真 | 是 | 多 | 么 | 只 | 是 | 到 | 了 | 这 | 时 | 候 |， | 我 | 的 | 往 | 事 | 才 | 在 | 我 | 眼 | 前 | 幻 | 现 | 得 | 清 | 晰 |， | 母 | 亲 | 的 | 苦 | 难 | 与 | 伟 | 大 | 才 | 在 | 我 | 心 | 中 | 摇 | 晃 | 着 | 轮 | 椅 | 在 | 园 | 中 | 慢 | 慢 | 走 |， | 又 | 是 | 雾 | 罩 | 的 | 清 | 晨 |， | 又 | 是 | 新 | 阳 | 高 | 悬 | 的 | 白 | 昼 |， | 我 | 只 | 想 | 着 | 一 | 件 | 事 |： | 曾 | 有 | 过 | 好 | 多 | 回 |， | 我 | 在 | 这 | 园 | 子 | 里 | 呆 | 得 | 太 | 久 | 了 |， | 母 | 亲 | 就 | 来 | 找 | 我 |。 | 她 | 来 | 找 | 我 | 又 | 不 | 想 | 让 | 我 | 发 | 觉 |

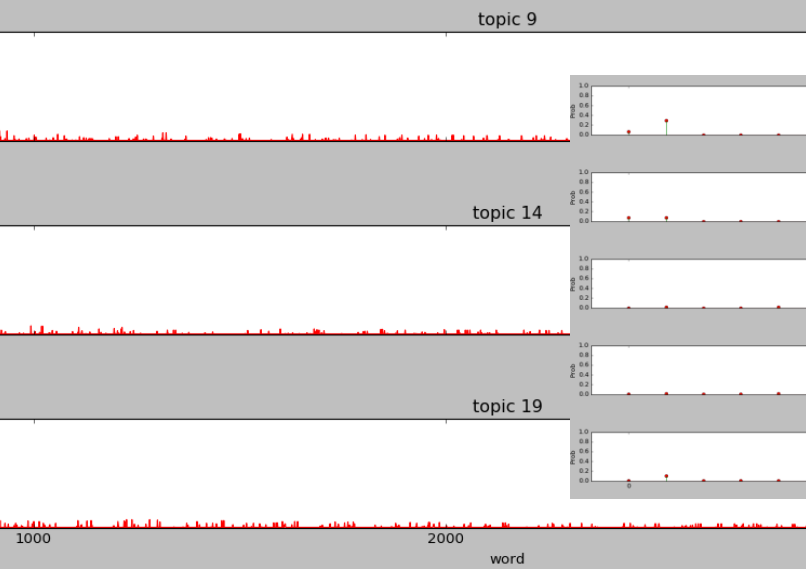
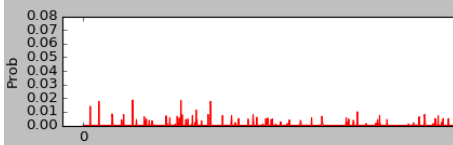
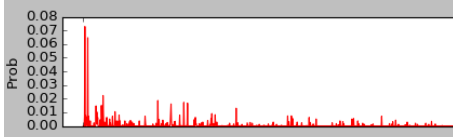
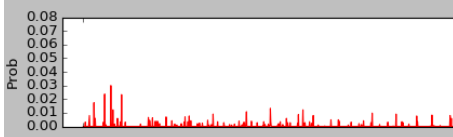
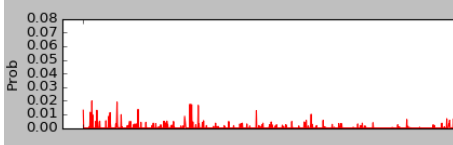
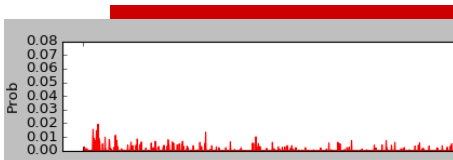


Jason Bell. *Machine Learning: Hands-On for Developers and Technical Professionals*. Wiley. 2015

# LDA

文档 1 : 茶馆 ( 0.0163591635916 ) 社会 ( 0.00528905289053 ) 王利发 ( 0.00528905289053 )  
文档 2 : 决议 ( 0.0138983050847 ) 打击 ( 0.00824858757062 ) 安理会 ( 0.00824858757062 )  
文档 3 : 会议 ( 0.0124491456469 ) 脱贫 ( 0.0124491456469 ) 党校 ( 0.0108218063466 )  
文档 4 : 美团 ( 0.0306066176471 ) 阿里 ( 0.0103860294118 ) 业务 ( 0.0103860294118 )  
文档 5 : 户口 ( 0.0221347331584 ) 登记 ( 0.0195100612423 ) 人口 ( 0.0142607174103 )  
文档 6 : 人员 ( 0.0111471861472 ) 飞机 ( 0.00898268398268 ) 称 ( 0.00681818181818 )  
文档 7 : 号线 ( 0.0328544061303 ) 站 ( 0.0194444444444 ) 14 ( 0.0184865900383 )  
文档 8 : 支付 ( 0.0198394495413 ) 腾讯 ( 0.0072247706422 ) 支付宝 ( 0.0072247706422 )  
文档 9 : 决议 ( 0.0138983050847 ) 打击 ( 0.00824858757062 ) 安理会 ( 0.00824858757062 )  
文档 10 : 足协 ( 0.0186473429952 ) 足球 ( 0.0138164251208 ) 佩兰 ( 0.011884057971 )

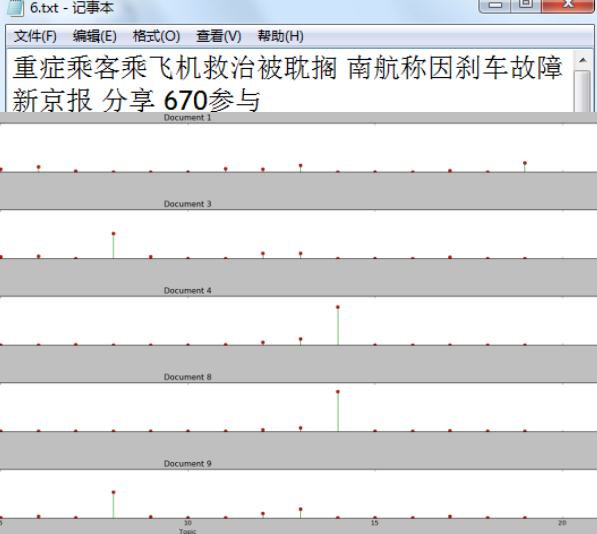
主题 1 : 美团 ( 0.0306066176471 ) 阿里 ( 0.0103860294118 ) 业务 ( 0.0103860294118 )  
主题 2 : 会议 ( 0.0124491456469 ) 脱贫 ( 0.0124491456469 ) 党校 ( 0.0108218063466 )  
主题 3 : 号线 ( 0.0328544061303 ) 站 ( 0.0194444444444 ) 14 ( 0.0184865900383 )  
主题 4 : 人物 ( 0.00214876033058 ) 民族 ( 0.00214876033058 ) 资本家 ( 0.00214876033058 )  
主题 5 : 足协 ( 0.0186473429952 ) 足球 ( 0.0138164251208 ) 佩兰 ( 0.011884057971 )  
主题 6 : 户口 ( 0.0221347331584 ) 登记 ( 0.0195100612423 ) 人口 ( 0.0142607174103 )  
主题 7 : 决议 ( 0.0138983050847 ) 打击 ( 0.00824858757062 ) 安理会 ( 0.00824858757062 )  
主题 8 : 人员 ( 0.0111471861472 ) 飞机 ( 0.00898268398268 ) 称 ( 0.00681818181818 )  
主题 9 : 茶馆 ( 0.0163591635916 ) 社会 ( 0.00528905289053 ) 王利发 ( 0.00528905289053 )  
主题 10 : 支付 ( 0.0198394495413 ) 腾讯 ( 0.0072247706422 ) 支付宝 ( 0.0072247706422 )



1.txt - 记事本  
文件(F) 编辑(E) 格式(O) 查看(V) 帮助(H)  
话剧《茶馆》是著名作家老舍先生创作的一部不朽的名著，三幕话剧剧本，1957年完成；1958年由北京人民艺术剧院首排，焦菊隐、夏淳导演，于是之、郑榕、蓝天野、英若诚、黄宗洛等人主演，全剧以老北京一家大茶馆的兴衰变迁为背景，向人们展示了从清末到抗战胜利后的50年间北京的社会风貌及各阶层人物的不同命运。

3.txt - 记事本  
文件(F) 编辑(E) 格式(O) 查看(V) 帮助(H)  
工作的意见》，听取关于巡视55家国有重要骨干企业有关情况的专题报告。中共中央总书记习近平主持会议。  
会议认为，改革开放以来，党和国家实施大规模扶贫开发，使7亿农村贫困人口摆脱贫困，取得了举世瞩目的伟大成就，谱写了人类反贫困历史上的辉煌篇章。党的十八大以来，我们把扶贫开发工作摆在更加突出的位置，实施精准扶贫，开创了扶贫开发事业新局面。

6.txt - 记事本  
文件(F) 编辑(E) 格式(O) 查看(V) 帮助(H)  
重症乘客乘飞机救治被耽搁 南航称因刹车故障  
新京报 分享 670参与



之后其腹痛的情况越来越严重，空乘人员赶紧帮忙预约了救护车，空乘和急救人员被指为谁该抬患者下飞机发生争执，患者最后自行勉强下旋梯爬进救护车。最后张先生在8小时后辗转了首都机场医院等，才被推送到北京大学人

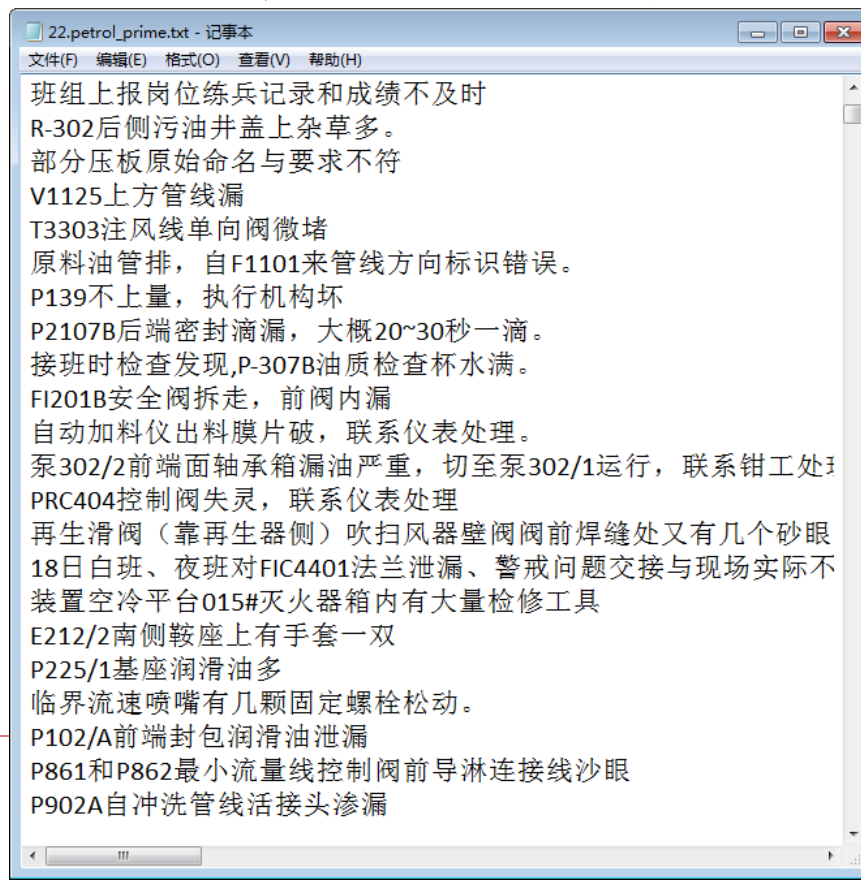


# 例：石油例检结果处理

□ 针对国内某石油企业的例行检查处理结果，试通过主题模型方案，分析例检结果中最突出的问题是什么？

■ 文本共4700个，

■ 单个文档十数字



# 其他内容

## □ 最大熵模型

- 自然语言处理解决标记问题

## □ 聚类

- K-means/K-Medoids/密度聚类/谱聚类

## □ 降维

- PCA/SVD/ICA

## □ SVM

- 与核技术相结合

## □ 主题模型pLSA/LDA

- 与聚类、标签传递算法相结合

## □ 条件随机场

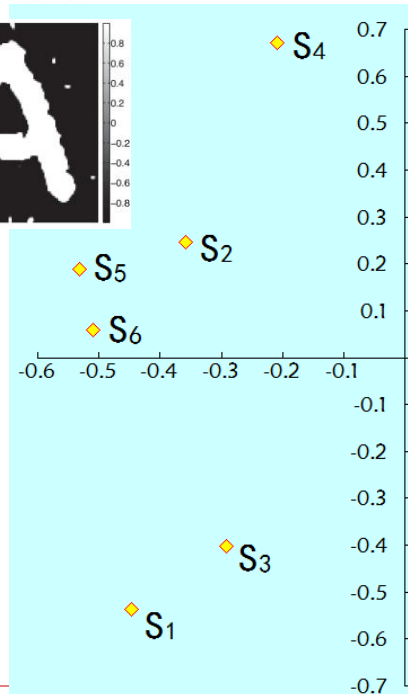
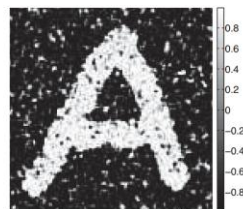
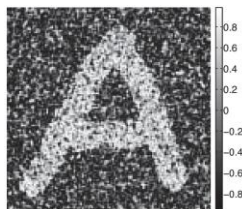
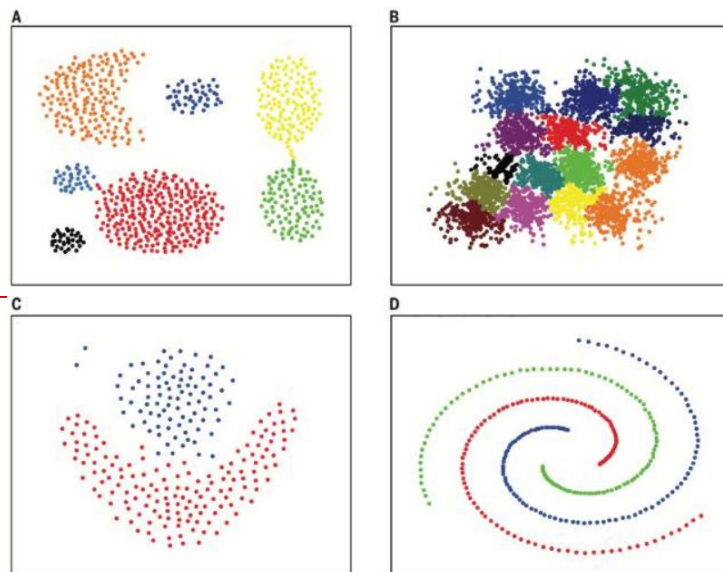
- 无向图模型，链式条件随机场解决标记问题

## □ 变分推导Variation Inference

- 与EM、贝叶斯相结合，参数、隐变量的学习

## □ 深度学习

- 大规模人工神经网络





# 本课程参考文献(部分)

---

- Christopher M. Bishop, Pattern Recognition and Machine Learning, Springer-Verlag, 2006
- Kevin P. Murphy, Machine Learning:A Probabilistic Perspective, The MIT Press, 2012
- 李航, 统计学习方法, 清华大学出版社, 2012
- Stephen Boyd,Lieven Vandenberghe, Convex Optimization, Cambridge University Press, 2004
- Thomas M. Cover, Joy A. Thomas, Elements of Information Theory,2006
- 各章节特定的经典论文, 如:
  - Alex Rodriguez, Alessandro Laio, Clustering by fast search and find of density peak, Science 344.6191(2014)
  - David M. Blei, Andrew Y. Ng, Michael I. Jordan, Latent Dirichlet Allocation, 2003

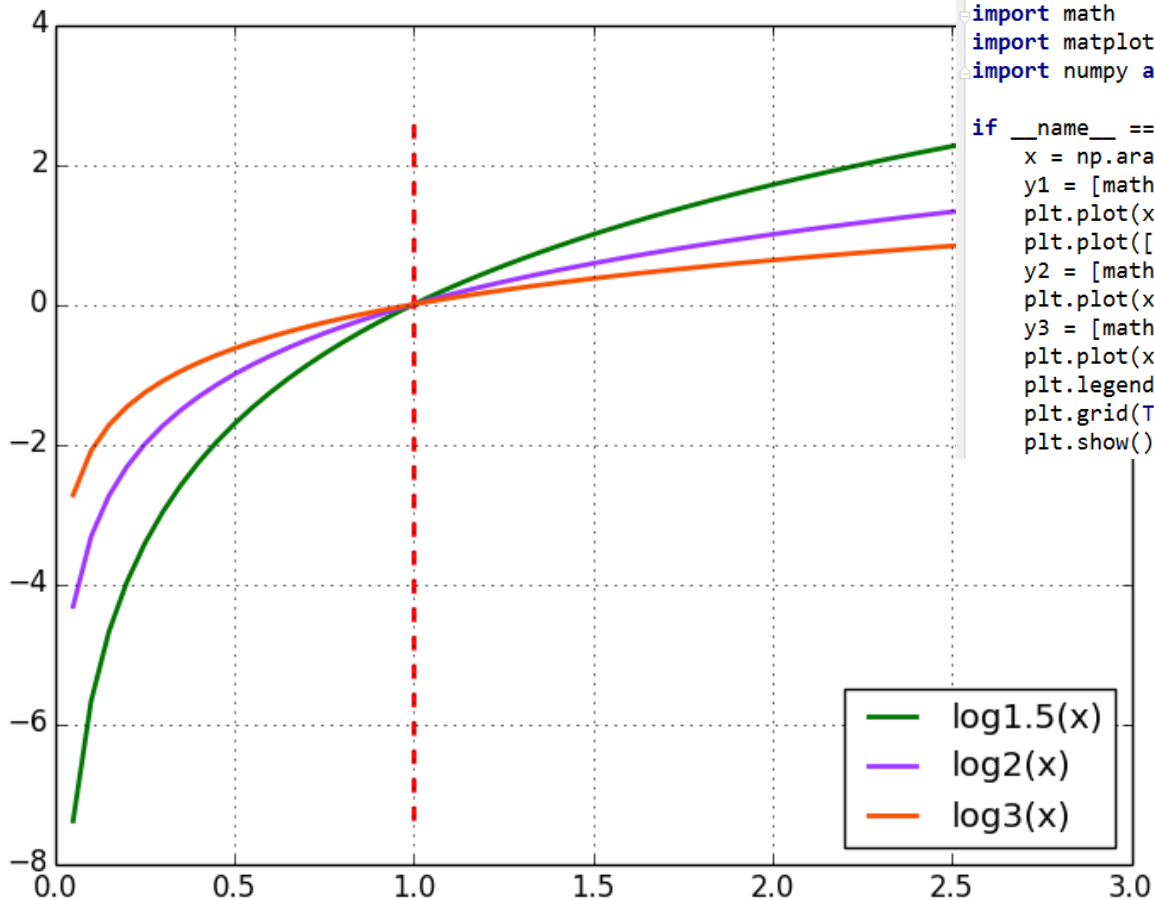
# 回忆知识

---

□ 求S的值：

$$S = \frac{1}{0!} + \frac{1}{1!} + \frac{1}{2!} + \frac{1}{3!} + \frac{1}{4!} + \dots + \frac{1}{n!} + \dots$$

# 对数函数的上升速度



```
import math
import matplotlib.pyplot as plt
import numpy as np

if __name__ == "__main__":
    x = np.arange(0.05, 3, 0.05)
    y1 = [math.log(a, 1.5) for a in x]
    plt.plot(x, y1, linewidth=2, color='#007500', label='log1.5(x)')
    plt.plot([1, 1], [y1[0], y1[-1]], "r--", linewidth=2)
    y2 = [math.log(a, 2) for a in x]
    plt.plot(x, y2, linewidth=2, color='#9F35FF', label='log2(x)')
    y3 = [math.log(a, 3) for a in x]
    plt.plot(x, y3, linewidth=2, color='#F75000', label='log3(x)')
    plt.legend(loc='lower right')
    plt.grid(True)
    plt.show()
```

# 问题分析

□ 令  $f(x) = \log_a x$

□ 则：

$$\frac{f(x + \Delta x) - f(x)}{\Delta x} = \frac{\log_a(x + \Delta x) - \log_a x}{\Delta x} = \frac{\log_a\left(\frac{x + \Delta x}{x}\right)}{\Delta x} = \log_a\left(\frac{x + \Delta x}{x}\right)^{\frac{1}{\Delta x}}$$

$$\xrightarrow{\because x=1} \log_a(1 + \Delta x)^{\frac{1}{\Delta x}} \stackrel{\text{令}}{=} 1 \Rightarrow \lim_{\Delta x \rightarrow 0} (1 + \Delta x)^{\frac{1}{\Delta x}} = a$$

□ 问： $\lim_{n \rightarrow \infty} \left(1 + \frac{1}{n}\right)^n = ?$

# 附：构造数列 $\{x_n\}$

$$\begin{aligned}x_n &= \left(1 + \frac{1}{n}\right)^n \\&= 1 + C_n^1 \frac{1}{n} + C_n^2 \frac{1}{n^2} + C_n^3 \frac{1}{n^3} + \cdots + C_n^n \frac{1}{n^n} \\&= 1 + n \cdot \frac{1}{n} + \frac{n(n-1)}{2!} \cdot \frac{1}{n^2} + \frac{n(n-1)(n-2)}{3!} \cdot \frac{1}{n^3} + \cdots + \frac{n(n-1)(n-2)\cdots 1}{n!} \cdot \frac{1}{n^n} \\&= 1 + 1 + \frac{1}{2!} \cdot \left(1 - \frac{1}{n}\right) + \frac{1}{3!} \cdot \left(1 - \frac{1}{n}\right) \left(1 - \frac{2}{n}\right) + \cdots + \frac{1}{n!} \cdot \left(1 - \frac{1}{n}\right) \left(1 - \frac{2}{n}\right) \cdots \left(1 - \frac{n-1}{n}\right) \\&< 1 + 1 + \frac{1}{2!} + \frac{1}{3!} + \cdots + \frac{1}{n!} \\&< 1 + 1 + \frac{1}{2} + \frac{1}{2^2} + \cdots + \frac{1}{2^{n-1}} \\&= 3 - \frac{1}{2^{n-1}} \\&< 3\end{aligned}$$

# 附：自然常数 $\lim_{x \rightarrow \infty} \left(1 + \frac{1}{x}\right)^x = e$

□ 根据前文中  $a_n = \left(1 + \frac{1}{n}\right)^n$  的二项展开式，已经证明数组  $\{a_n\}$  单增有上界，因此，必有极限，记做  $e$ 。

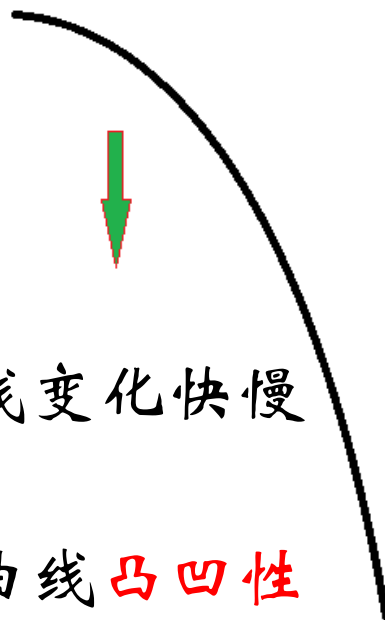
□ 同时：
$$\left(1 + \frac{1}{n+1}\right)^n < \left(1 + \frac{1}{x}\right)^x < \left(1 + \frac{1}{n}\right)^{n+1}$$

$$\lim_{n \rightarrow \infty} \left(1 + \frac{1}{n+1}\right)^n = \lim_{n \rightarrow \infty} \frac{\left(1 + \frac{1}{n+1}\right)^{n+1}}{1 + \frac{1}{n+1}} = \frac{\lim_{n \rightarrow \infty} \left(1 + \frac{1}{n+1}\right)^{n+1}}{\lim_{n \rightarrow \infty} \left(1 + \frac{1}{n+1}\right)} = \frac{e}{1+0} = e$$

$$\lim_{n \rightarrow \infty} \left(1 + \frac{1}{n}\right)^{n+1} = \lim_{n \rightarrow \infty} \left( \left(1 + \frac{1}{n}\right)^n \left(1 + \frac{1}{n}\right) \right) = \lim_{n \rightarrow \infty} \left(1 + \frac{1}{n}\right)^n \cdot \lim_{n \rightarrow \infty} \left(1 + \frac{1}{n}\right) = e \cdot (1+0) = e$$

□ 根据两边夹定理，函数  $f(x) = \left(1 + \frac{1}{x}\right)^x$  的极限存在，为  $e$ 。

# 导数



- 简单的说，导数就是曲线的斜率，是曲线变化快慢的反应
- **二阶导数**是斜率变化快慢的反应，表征曲线**凸凹性**
  - 二阶导数连续的曲线，往往称之为“**光顺**”的。
  - 还记得高中物理老师时常念叨的吗：**加速度**的方向总是指向轨迹曲线凹的一侧。
- 根据  $\lim_{x \rightarrow \infty} \left(1 + \frac{1}{x}\right)^x = e$  可以得到函数  $f(x) = \ln x$  的导数，进一步根据换底公式、反函数求导等，得到其他初等函数的导数。

# 常用函数的导数

---

$$C' = 0$$

$$(x^n)' = nx^{n-1}$$

$$(\sin x)' = \cos x$$

$$(\cos x)' = -\sin x$$

$$(a^x)' = a^x \ln a$$

$$(e^x)' = e^x$$

$$(\log_a x)' = \frac{1}{x} \log_a e$$

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

$$(u + v)' = u' + v'$$

$$(uv)' = u'v + uv'$$



# 应用1

---

- 已知函数  $f(x) = x^x, x > 0$
- 求  $f(x)$  的最小值
  - 领会幂指数函数的一般处理套路

- 附：
$$N^{\frac{1}{\log_2 N}} = ?$$

- 在计算机算法跳跃表Skip List的分析中，用到了该常数。
- 背景：跳表是支持增删改查的动态数据结构，能够达到与平衡二叉树、红黑树近似的效率，而代码实现简单。

# 求解 $x^x$

$$t = x^x$$

$$\rightarrow \ln t = x \ln x$$

$$\xrightarrow{\text{两边对 } x \text{ 求导}} \frac{1}{t} t' = \ln x + 1$$

$$\xrightarrow{\text{令 } t'=0} \ln x + 1 = 0$$

$$\rightarrow x = e^{-1}$$

$$\rightarrow t = e^{-\frac{1}{e}}$$

## 积分应用2: $N \rightarrow \infty \Rightarrow \ln N! \rightarrow N(\ln N - 1)$

$$\begin{aligned}\ln N! &= \sum_{i=1}^N \ln i \approx \int_1^N \ln x dx \\ &= x \ln x \Big|_1^N - \int_1^N x d \ln x \\ &= N \ln N - \int_1^N x \cdot \frac{1}{x} dx \\ &= N \ln N - x \Big|_1^N \\ &= N \ln N - N + 1 \\ &\rightarrow N \ln N - N\end{aligned}$$

# Taylor公式 – Maclaurin公式

---

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

$$f(x) = f(0) + f'(0)x + \frac{f''(0)}{2!}x^2 + \dots + \frac{f^{(n)}(0)}{n!}x^n + o(x^n)$$

# Taylor公式的应用1

□ 数值计算：初等函数值的计算(在原点展开)

$$\sin x = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} + \frac{x^9}{9!} + \cdots + (-1)^{m-1} \frac{x^{2m-1}}{(2m-1)!} + R_{2m}$$

$$e^x = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \cdots + \frac{x^n}{n!} + R_n$$

□ 在实践中，往往需要做一定程度的变换。

# Taylor公式的应用1：计算 $e^x$

□ 给定正实数 $x$ ，计算 $e^x=?$

□ 一种可行的思路：

□ 求整数 $k$ 和小数 $r$ ，使得

■  $x = k \cdot \ln 2 + r, |r| \leq 0.5 \cdot \ln 2$

□ 从而：
$$e^x = e^{k \cdot \ln 2 + r}$$

$$= e^{k \cdot \ln 2} \cdot e^r$$

$$= 2^k \cdot e^r$$

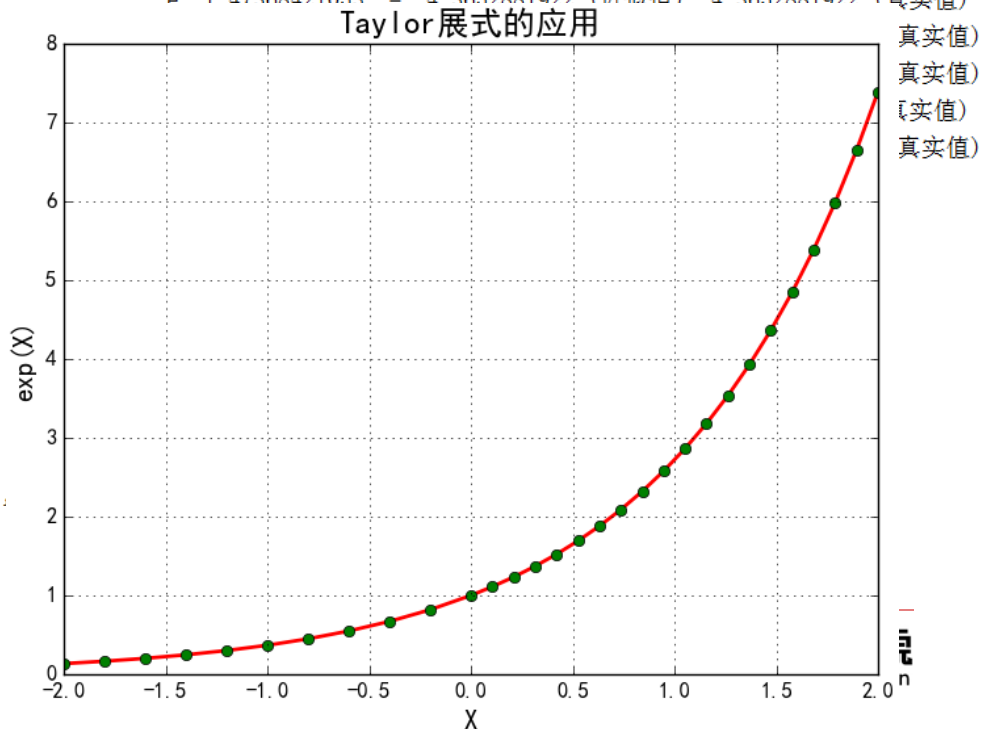
# Taylor展式的应用

```
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)



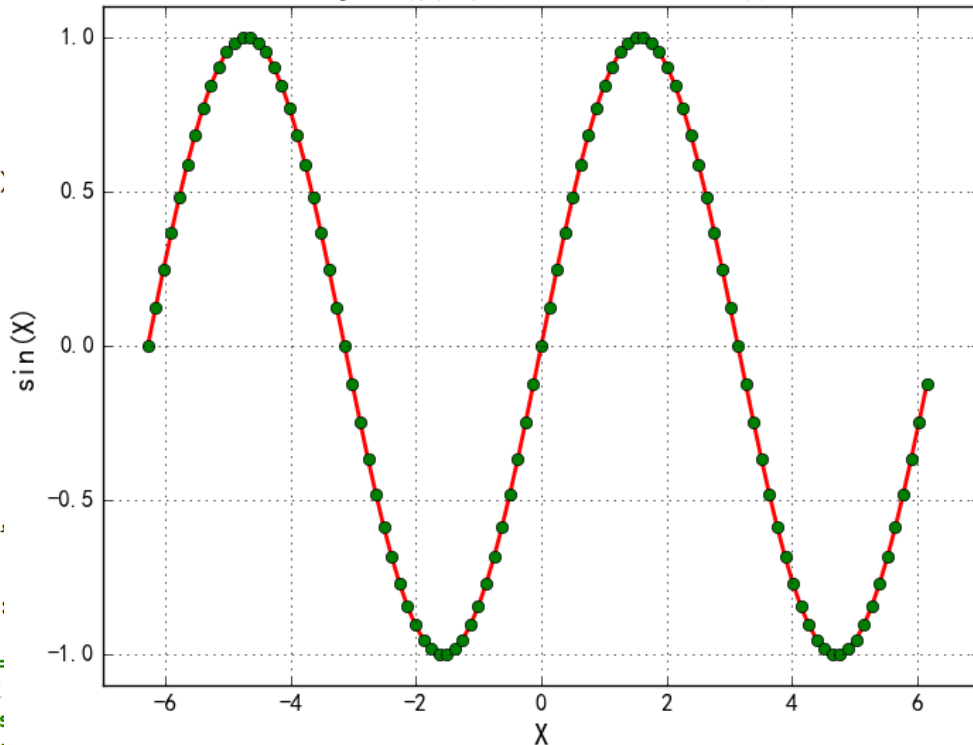
# Taylor展式的应用

```
def calc_sin_small(x):  
    x2 = -x ** 2  
    t = x  
    f = 1  
    sum = 0  
    for i in range(10):  
        sum += t / f  
        t *= x2  
        f *= ((2*i+2)*(2*i+3))  
    return sum
```

```
def calc_sin(x):  
    a = x / (2*np.pi)  
    k = np.floor(a)  
    a = x - k*2*np.pi  
    return calc_sin_small(a)
```

```
if __name__ == "__main__":  
    t = np.linspace(-2*np.pi,  
                    # 横轴数据  
                    y = np.empty_like(t)  
    for i, x in enumerate(t):  
        y[i] = calc_sin(x)  
        print 'sin(', x, ') :  
        # print '误差: ', y[i]  
    mpl.rcParams['font.sans-serif'] = ['SimHei']  
    mpl.rcParams['axes.unicode_minus'] = False  
    plt.figure(facecolor='w')  
    plt.plot(t, y, 'r-', t, y, 'go', linewidth=2)  
    plt.title(u'Taylor展式的应用 - 正弦函数', fontsize=18)  
    plt.xlabel('X', fontsize=15)  
    plt.ylabel('sin(X)', fontsize=15)  
    plt.xlim((-7, 7))  
    plt.ylim((-1.1, 1.1))  
    plt.grid(True)  
    plt.show()
```

Taylor展式的应用 - 正弦函数



sin(x)	(近似值)	(真实值)
sin(-1.25663706144)	-0.951066447544	-0.951056516295
sin(-1.13097335529)	-0.904843692145	-0.904827052466
sin(-1.00530964915)	-0.844355457307	-0.844327925502
sin(-0.879645943005)	-0.770558254809	-0.770513242776
sin(-0.753982236862)	-0.684619861245	-0.684547105929
sin(-0.628318530718)	-0.587901575106	-0.587785252292
sin(-0.502654824574)	-0.481937724358	-0.481753674102
sin(-0.376991118431)	-0.368412871668	-0.368124552685
sin(-0.251327412287)	-0.249137247033	-0.248689887165

1027336 (近似值)	-0.125333233564 (真实值)	
197e-16 (近似值)	8.881784197e-16 (真实值)	
33564 (近似值)	0.125333233564 (真实值)	
87165 (近似值)	0.248689887165 (真实值)	
52685 (近似值)	0.368124552685 (真实值)	
74102 (近似值)	0.481753674102 (真实值)	
52292 (近似值)	0.587785252292 (真实值)	
05929 (近似值)	0.684547105929 (真实值)	
42776 (近似值)	0.770513242776 (真实值)	
5502 (近似值)	0.844327925502 (真实值)	
2466 (近似值)	0.904827052466 (真实值)	
6295 (近似值)	0.951056516295 (真实值)	
0729 (近似值)	0.982287250729 (真实值)	
8428 (近似值)	0.998026728428 (真实值)	
8428 (近似值)	0.998026728428 (真实值)	
0729 (近似值)	0.982287250729 (真实值)	
6295 (近似值)	0.951056516295 (真实值)	
466 (近似值)	0.904827052466 (真实值)	
5502 (近似值)	0.844327925502 (真实值)	
2775 (近似值)	0.770513242776 (真实值)	
sin( 2.38701041073 )	= 0.684547105927 (近似值)	0.684547105929 (真实值)
sin( 2.51327412287 )	= 0.587785252288 (近似值)	0.587785252292 (真实值)
sin( 2.63893782902 )	= 0.481753674088 (近似值)	0.481753674102 (真实值)
sin( 2.76460153516 )	= 0.368124552648 (近似值)	0.368124552685 (真实值)

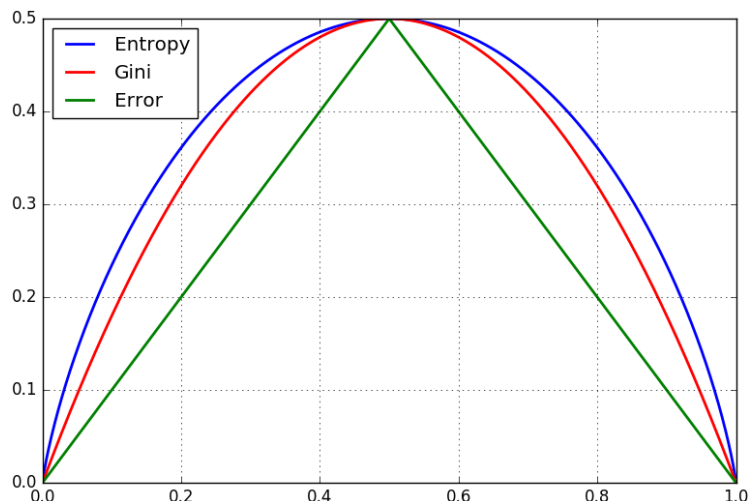


# Taylor公式的应用2

□ 考察Gini系数的图像、熵、分类误差率三者之间的关系

■ 将 $f(x)=-\ln x$ 在 $x=1$ 处一阶展开，忽略高阶无穷小，得到 $f(x)\approx 1-x$

$$H(X) = -\sum_{k=1}^K p_k \ln p_k$$
$$\approx \sum_{k=1}^K p_k (1 - p_k)$$

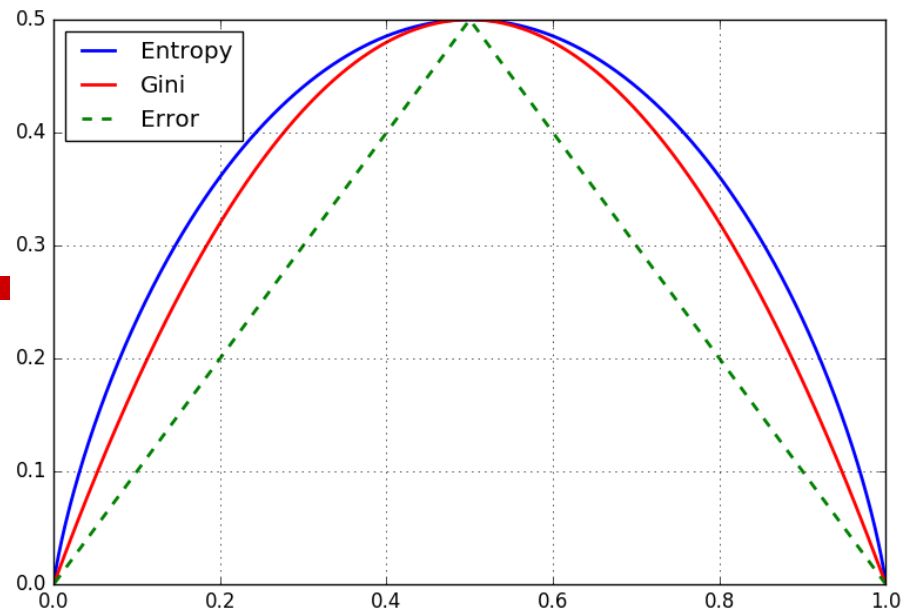


■ 上述结论，在决策树章节中会进一步讨论

# Gini系数的生成

```
import numpy as np
import matplotlib.pyplot as plt

if __name__ == "__main__":
    p = np.arange(0.001, 1, 0.001, dtype=np.float)
    gini = 2 * p * (1-p)
    h = -(p * np.log2(p) + (1-p) * np.log2(1-p))/2
    err = 1 - np.max(np.vstack((p, 1-p)), 0)
    plt.plot(p, h, 'b-', linewidth=2, label='Entropy')
    plt.plot(p, gini, 'r-', linewidth=2, label='Gini')
    plt.plot(p, err, 'g-', linewidth=2, label='Error')
    plt.grid(True)
    plt.legend(loc='upper left')
    plt.show()
```



# 方向导数

- 如果函数  $z=f(x,y)$  在点  $P(x,y)$  是可微分的，那么，函数在该点沿任一方向  $L$  的方向导数都存在，且有：

$$\frac{\partial f}{\partial l} = \frac{\partial f}{\partial x} \cos \varphi + \frac{\partial f}{\partial y} \sin \varphi$$

- 其中， $\psi$  为  $x$  轴到方向  $L$  的转角。

# 梯度

- 设函数 $z=f(x,y)$ 在平面区域 $D$ 内具有一阶连续偏导数, 则对于每一个点 $P(x,y) \in D$ , 向量

$$\left( \frac{\partial f}{\partial x}, \frac{\partial f}{\partial y} \right)$$

为函数 $z=f(x,y)$ 在点 $P$ 的梯度, 记做 $\text{grad}f(x,y)$

- 梯度的方向是函数在该点变化最快的方向
  - 考虑一座解析式为 $z=H(x,y)$ 的山, 在 $(x_0,y_0)$ 的梯度是在该点坡度变化最快的方向。
- 梯度下降法
  - 思考: 若下山方向和梯度呈 $\theta$ 夹角, 下降速度是多少?

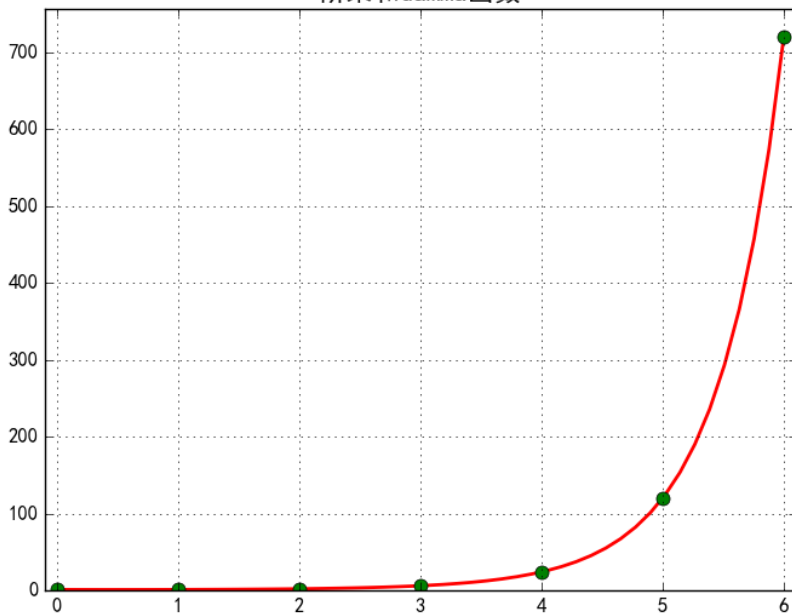
# Γ函数

$$\Gamma(x) = (x-1) \cdot \Gamma(x-1) \Rightarrow \frac{\Gamma(x)}{\Gamma(x-1)} = x-1$$

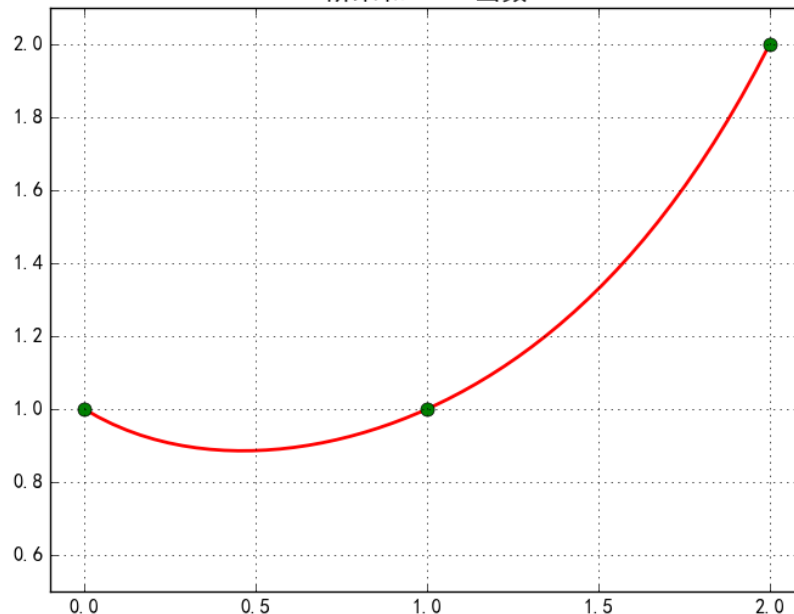
□ Γ函数是阶乘在实数上的推广

$$\Gamma(x) = \int_0^{+\infty} t^{x-1} e^{-t} dt = (x-1)!$$

阶乘和Gamma函数



阶乘和Gamma函数

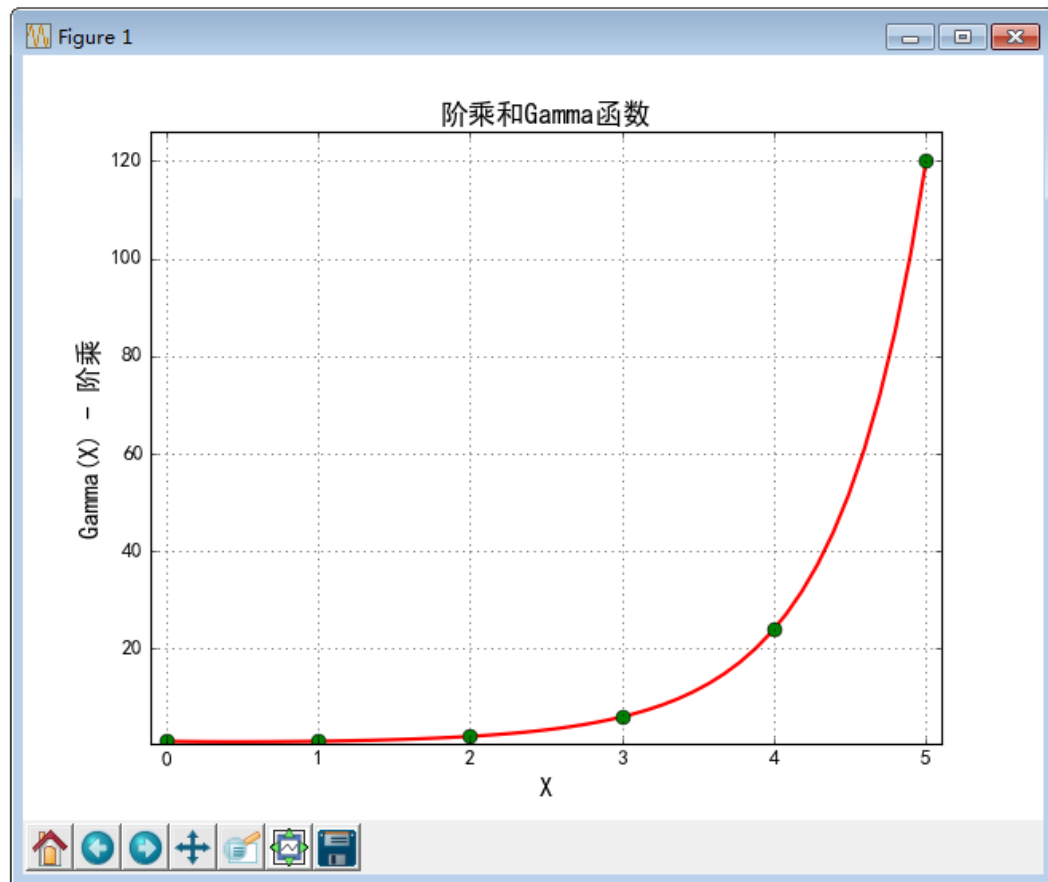


# Code

```
import numpy as np
import matplotlib as mpl
import matplotlib.pyplot as plt
from scipy.special import gamma
from scipy.special import factorial

mpl.rcParams['axes.unicode_minus'] = False
mpl.rcParams['font.sans-serif'] = 'SimHei'

if __name__ == '__main__':
    N = 5
    x = np.linspace(0, N, 50)
    y = gamma(x+1)
    plt.figure(facecolor='w')
    plt.plot(x, y, 'r-', lw=2)
    z = np.arange(0, N+1)
    f = factorial(z, exact=True)    # 阶乘
    print f
    plt.plot(z, f, 'go', markersize=8)
    plt.grid(b=True)
    plt.xlim(-0.1, N+0.1)
    plt.ylim(0.5, np.max(y)*1.05)
    plt.xlabel(u'X', fontsize=15)
    plt.ylabel(u'Gamma(X) - 阶乘', fontsize=15)
    plt.title(u'阶乘和Gamma函数', fontsize=16)
    plt.show()
```



5.7.gamma

C:\Python27\python.exe D:/Python/ML/5.Python/5.7.gamma.py

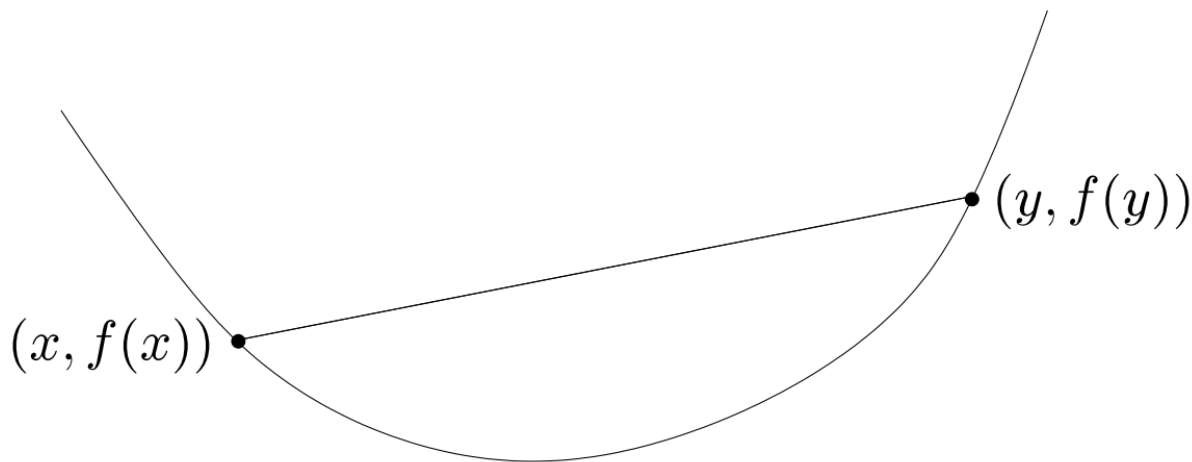
[ 1 1 2 6 24 120]

# 凸函数

□ 若函数 $f$ 的定义域 $\text{dom}f$ 为凸集，且满足

$\forall x, y \in \text{dom} f, 0 \leq \theta \leq 1$ ，有

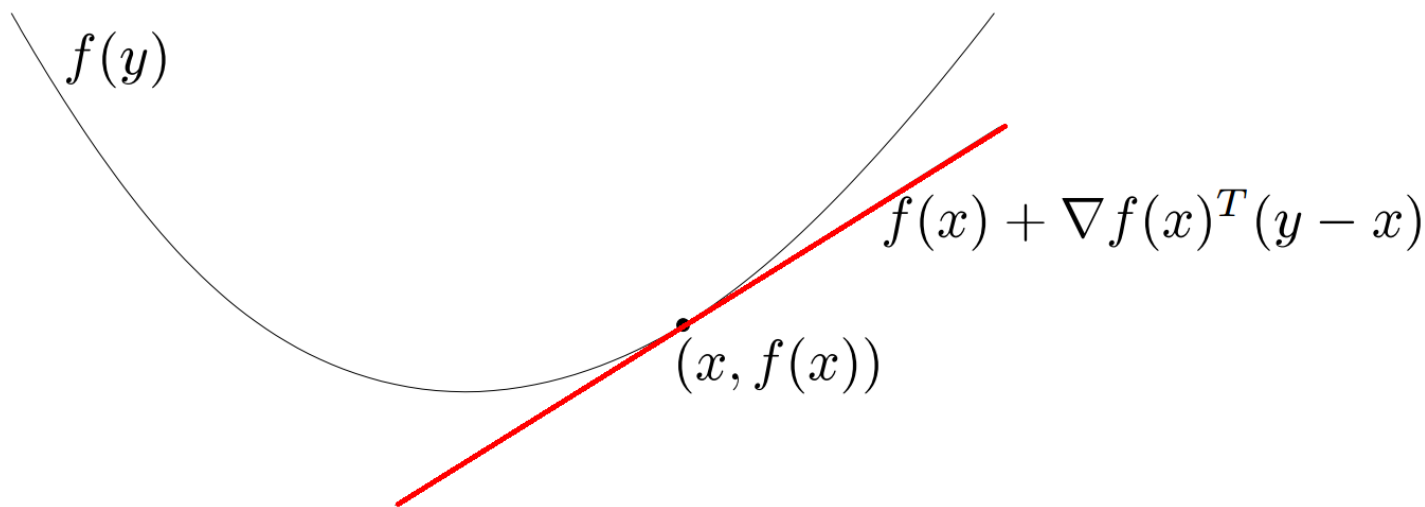
$$f(\theta x + (1 - \theta)y) \leq \theta f(x) + (1 - \theta)f(y)$$



# 一阶可微

□ 若 $f$ 一阶可微，则函数 $f$ 为凸函数当前仅当 $f$ 的定义域 $\text{dom}f$ 为凸集，且

$$\forall x, y \in \text{dom}f, f(y) \geq f(x) + \nabla f(x)^T (y - x)$$





# 二阶可微

---

- 若函数 $f$ 二阶可微，则函数 $f$ 为凸函数当前仅当 $\text{dom}$ 为凸集，且

$$\nabla^2 f(x) \succeq 0$$

- 若 $f$ 是一元函数，上式表示二阶导大于等于0
- 若 $f$ 是多元函数，上式表示二阶导Hessian矩阵半正定。

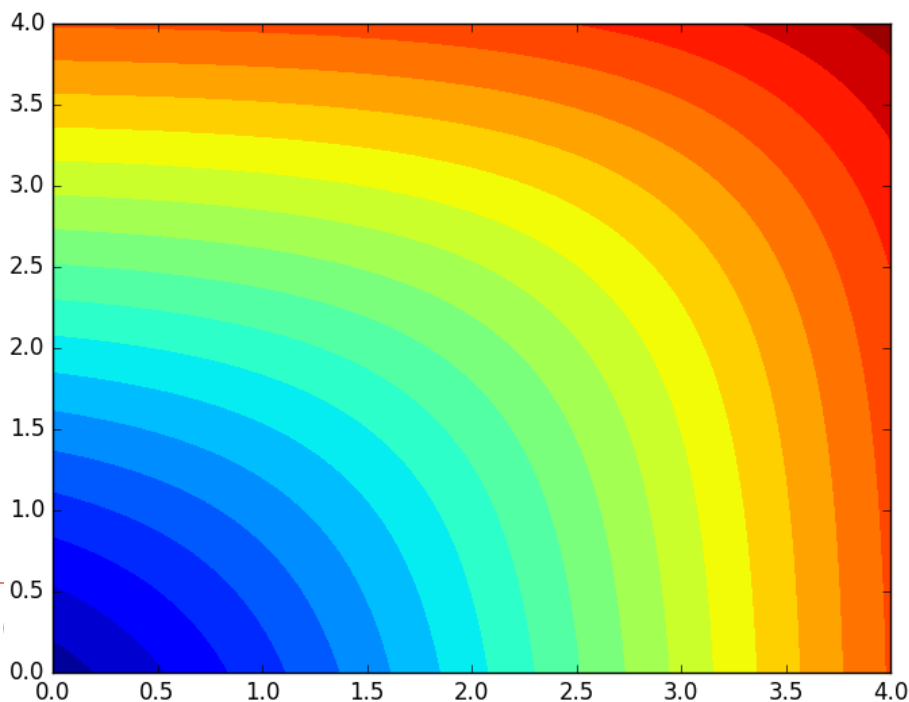
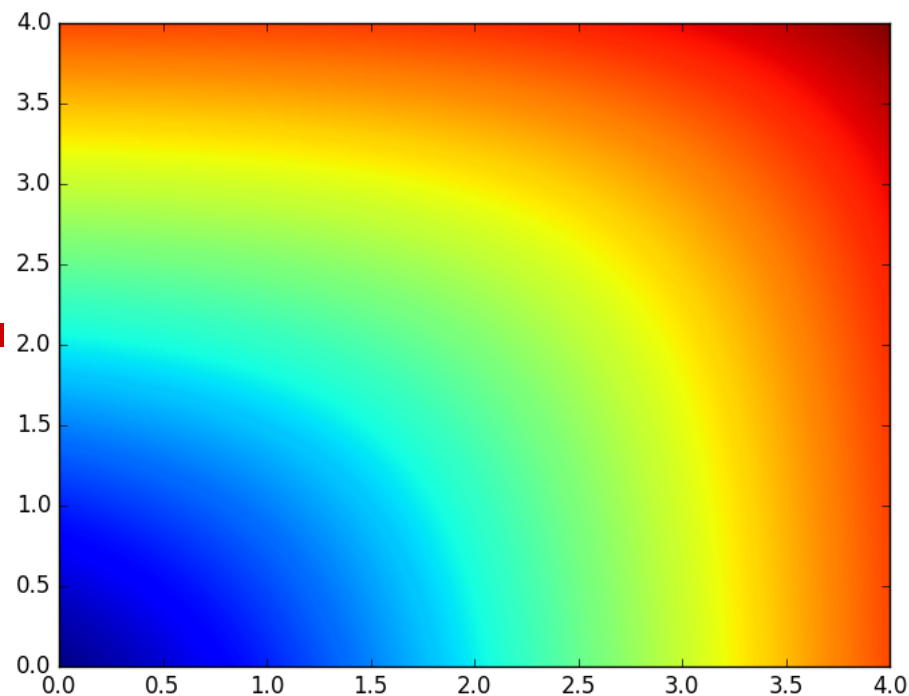
# 凸函数举例

- 指数函数  $f(x) = e^{ax}$
- 幂函数  $f(x) = x^a$ ,  $x \in R^+$ ,  $a \geq 1$  或  $a \leq 0$
- 负对数函数  $f(x) = -\ln x$
- 负熵函数  $f(x) = x \ln x$
- 范数函数  $f(\vec{x}) = \|\vec{x}\|$
- 最大值函数  $f(\vec{x}) = \max(x_1, x_2, \dots, x_n)$
- 指数线性函数  $f(\vec{x}) = \log(e^{x_1} + e^{x_2} + \dots + e^{x_n})$ 
  - Log Sum Exp  $\rightarrow \max(x_1, x_2, \dots, x_n)$

# Code

$$f(x_1, x_2) = \log(e^{x_1} + e^{x_2})$$

```
if __name__ == "__main__":  
    fig = plt.figure()  
    ax = fig.add_subplot(111)  
    u = np.linspace(0, 4, 1000)  
    x, y = np.meshgrid(u, u)  
    z = np.log(np.exp(x) + np.exp(y))  
    ax.contourf(x, y, z, 20)  
    plt.show()
```



# 概率论

- 对概率的认识:  $P(x) \in [0,1]$ 
  - $P=0$ : 事件出现的概率为0 → 事件不会发生?
  - 若 $x$ 为离散/连续变量, 则 $P(x=x_0)$ 表示 $x_0$ 发生的概率/概率密度
- 累计分布函数:  $\Phi(x)=P(x \leq x_0)$ 
  - $\Phi(x)$ 一定为单增函数
  - $\min(\Phi(x))=0, \max(\Phi(x))=1$
- 思考: 将值域为 $[0,1]$ 的某单增函数 $y=F(x)$ 看成 $X$ 事件的累积概率函数
  - 若 $y=F(x)$ 可导, 则 $f(x)=F'(x)$ 为某概率密度函数
- P.S.
  - cumulative distribution function, CDF/Probability Density Function, pdf

# 古典概型

---

- 举例：将 $n$ 个不同的球放入 $N(N \geq n)$ 个盒子中，假设盒子容量无限，求事件 $A = \{\text{每个盒子至多有1个球}\}$ 的概率。

解 
$$P(A) = \frac{P_N^n}{N^n}$$

---

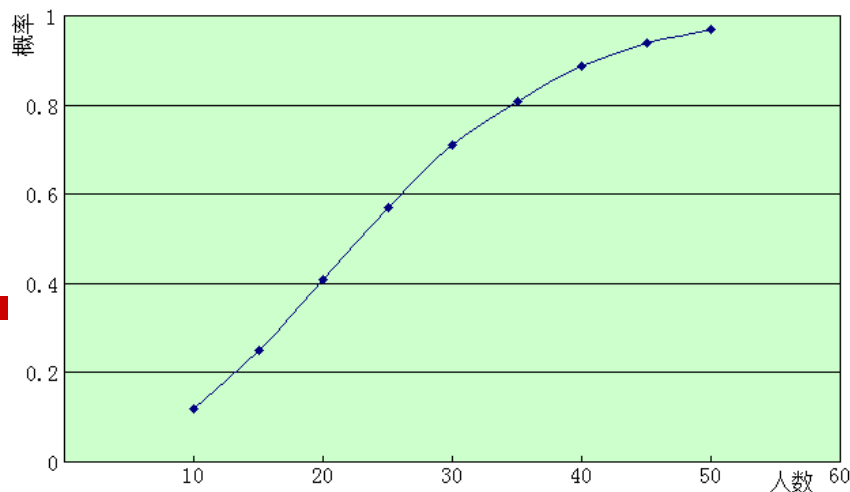
□ 基本事件总数：

- 第1个球，有N种放法；
- 第2个球，有N种放法；
- .....
- 共： $N^n$ 种放法。

□ 每个盒子至多放1个球的事件数：

- 第1个球，有N种放法；
- 第2个球，有N-1种放法；
- 第3个球，有N-2种放法；
- .....
- 共： $N(N-1)(N-2)\cdots(N-n+1) = P_N^n$

# 生日悖论



□ 假定小象学院《机器学习》班有50位同学，则至少有2人生日相同的概率是多少？

n	10	15	20	25	30	35	40	45	50
P	0.12	0.25	0.41	0.57	0.71	0.81	0.89	0.94	0.97

# 装箱问题

---

- 将12件正品和3件次品随机装在3个箱子中，每箱装5件，则每箱中恰有1件次品的概率是多少？



# 解

---

- 将15件产品装入3个箱子，每箱装5件，共有  $15!/(5!5!5!)$  种装法；
- 先把3件次品放入3个箱子，有  $3!$  种装法。对于这样的每一种装法，把其余12件产品装入3个箱子，每箱装4件，共有  $12!/(4!4!4!)$  种装法；
- $P(A) = (3! * 12! / (4!4!4!)) / (15! / (5!5!5!)) = 25/91$

# 与组合数的关系

- 把n个物品分成k组，使得每组物品的个数分别为 $n_1, n_2, \dots, n_k$ ，( $n = n_1 + n_2 + \dots + n_k$ )，则不同的分组方法有  $\frac{n!}{n_1! n_2! n_3! \dots n_k!}$  种。
- 上述问题的简化版本，即n个物品分成2组，第一组m个，第二组n-m个，则分组方法有  $\frac{n!}{m!(n-m)!}$ ，即： $C_n^m$ 。

# 组合数背后的秘密 $N \rightarrow \infty \Rightarrow \ln N! \rightarrow N(\ln N - 1)$

$$\begin{aligned} H &= \frac{1}{N} \ln \frac{N!}{\prod_{i=1}^k n_i!} = \frac{1}{N} \ln(N!) - \frac{1}{N} \sum_{i=1}^k \ln(n_i!) \\ &\rightarrow (\ln N - 1) - \frac{1}{N} \sum_{i=1}^k n_i (\ln n_i - 1) \\ &= \ln N - \frac{1}{N} \sum_{i=1}^k n_i \ln n_i = -\frac{1}{N} \left( \left( \sum_{i=1}^k n_i \ln n_i \right) - N \ln N \right) \\ &= -\frac{1}{N} \sum_{i=1}^k (n_i \ln n_i - n_i \ln N) = -\frac{1}{N} \sum_{i=1}^k \left( n_i \ln \frac{n_i}{N} \right) \\ &= -\sum_{i=1}^k \left( \frac{n_i}{N} \ln \frac{n_i}{N} \right) \rightarrow -\sum_{i=1}^k (p_i \ln p_i) \end{aligned}$$

# 参考文献

---

- Prof. Andrew Ng, *Machine Learning*, Stanford University
- 同济大学数学教研室, 高等数学, 高等教育出版社, 1996
- 王松桂、程维虎、高旅端, 概率论与数理统计, 科学出版社, 2000

# 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

■ 大数据分析挖掘

The screenshot shows the website [wenda.chinahadoop.cn/explore/](http://wenda.chinahadoop.cn/explore/). The page features a navigation bar with '发现' (Discover) highlighted in a red circle. Below the navigation bar, there are tabs for '全部', '招聘求职', '机器学习', '大数据平台技术', 'DCon', '大数据行业应用', 'NoSQL数据库', '数据科学', and '江湖救急'. The main content area displays a list of questions and answers, including:

- Question: "yarn运行时一直重复这个info...好像没找到资源, 应该从哪里检查呢?" (yarn runtime always repeats this info... seems to not find resources, where should I check?)
- Question: "两种不同的相关推荐列表" (Two different recommendation lists)
- Question: "如何在Linux下配java的JDK?" (How to configure java's JDK under Linux?)
- Question: "sqoop把mysql数据导入Hbase报如图错误" (sqoop imports mysql data into Hbase with the error shown in the figure)
- Question: "泛化误差公式推导" (Generalization error formula derivation)
- Question: "kafkaOffsetMonitor打开页面以后无法显示内容?" (After opening the kafkaOffsetMonitor page, the content cannot be displayed?)
- Question: "markdown公式编辑\$符号不起作用" (Markdown formula editing \$ symbol does not work)
- Question: "hadoop-2.6.2-src源码编译成功之后找不到native下图一所示文件, 执行图三所示搜索命令也没有找到, 进入源码编译之后的目录如图二! 这个文件找不到怎么解决呢? 是编译没产生?" (After successful compilation of hadoop-2.6.2-src source code, the native file shown in Figure 1 cannot be found, and the search command shown in Figure 3 also cannot find it, and the directory after compilation is shown in Figure 2! How to solve this problem if the file cannot be found? Is it not generated during compilation?)
- Question: "opentsdb安装时出现72个warning, 是正常的么?" (When installing opentsdb, 72 warnings appear, is it normal?)
- Question: "关于在线广告和个性化推荐区别的一点浅见" (A slight insight into the difference between online advertising and personalized recommendation)

The right sidebar contains sections for '招聘求职', '大数据行业应用', '数据科学', '系统与编程', '云计算技术', '热门话题', '热门用户', and '机器学习'.

---

感谢大家！

恳请大家批评指正！