

说明

- 对于所有程序，是在python3.8实现的，都可以运行
- 如果提示缺少对应的第三方库，那么需要你先去pip3 install一下，可以加入网上一些下载比较快的镜像源。我建议是在配置pycharm里面的python解释器时，直接配置Anaconda的python环境，里面很多必要的库都有，比如sklearn机器学习库等，就可以不用去自行pip安装了，很方便
- 代码在复制的时候，请一定要注意缩进问题

规划

线性规划

scipy求解

- 需要知道目标函数(一般是求最大或者最小值)和约束条件
求解前转化为下面的标准形式

$$\begin{aligned} & \min c^T x \\ & s.t. \begin{cases} Ax \leq b \\ Aeq * x = beq \\ lb \leq x \leq ub \end{cases} \end{aligned}$$

```
1 from scipy import optimize
2 import numpy as np
3
4 # 求解函数
5 res = optimize.linprog(c, A, b, Aeq, beq, LB, UB, X0, OPTIONS)
6 # 目标函数最小值
7 print(res.fun)
8 # 最优解
9 print(res.x)
```

- 标准形式是 \leq ，如果是 \geq ，则在两边加上符号-

举例1

样例1： 求解下列线性规划问题

$$\begin{aligned} \max z &= 2x_1 + 3x_2 - 5x_3 \\ \text{s.t.} \quad &\begin{cases} x_1 + x_2 + x_3 = 7 \\ 2x_1 - 5x_2 + x_3 \geq 10 \\ x_1 + 3x_2 + x_3 \leq 12 \\ x_1, x_2, x_3 \geq 0 \end{cases} \end{aligned}$$

- 使用scipy求解z的最大值
 - c是目标函数的系数矩阵
 - A是化成标准的 \leq 式子的左边的系数矩阵
 - B是化成标准的 \leq 式子的右边的数值矩阵
 - Aeq是所有 $=$ 左边的系数矩阵，记得里面是 $[[[]]]$ 二维
 - Beq是所有 $=$ 右边的数值矩阵
 - 下面第11行-c加-，是因为此题求的是最大值，但是标准格式是求最小值，所以加负号
 - 另外上面的3个变量都大于0，这里可以使用 `bounds=(0, None)` ,`bounds=(min,max)`是范围，None代表无穷，如果不写bounds，那默认是 $(0, \text{None})$

```
1 | res = optimize.linprog(-c, A, B, Aeq, Beq, bounds=(0, None))
```

```
1 | from scipy import optimize
2 | import numpy as np
3 |
4 | # 确定c,A,b,Aeq,beq
5 | c = np.array([2, 3, -5])
6 | A = np.array([[2, 5, -1], [1, 3, 1]])
7 | B = np.array([-10, 12])
8 | Aeq = np.array([[1, 1, 1]])
9 | Beq = np.array([7])
10 | # 求解
11 | res = optimize.linprog(-c, A, B, Aeq, Beq)
12 | print(res)
13 |
```

- res

```

1      con: array([1.80713222e-09])
2      fun: -14.57142856564506
3      message: 'Optimization terminated successfully.'
4      nit: 5
5      slack: array([-2.24583019e-10,  3.85714286e+00])
6      status: 0
7      success: True
8      x: array([6.42857143e+00, 5.71428571e-01, 2.35900788e-10])
9
10     Process finished with exit code 0

```

- fun是目标函数最小值
- x是最优解，即上面的x1,x2,x3的最优解

举例2

求解下列线性规划问题

$$\begin{aligned}
 \max z &= 2x_1 + 3x_2 + x_3 \\
 \text{s.t.} \quad &\begin{cases} x_1 + 2x_2 + 4x_3 = 101 \\ x_1 + 4x_2 + 2x_3 \geq 8 \\ 3x_1 + 2x_2 \geq 6 \\ x_1, x_2, x_3 \geq 0 \end{cases}
 \end{aligned}$$

```

1  from scipy import optimize
2  import numpy as np
3
4  c = np.array([2, 3, 1])
5  A = np.array([[ -1,  -4,  -2], [ -3,  -2,  0]])
6  B = np.array([ -8,  -6])
7  Aeq = np.array([[1, 2, 4]])
8  Beq = np.array([101])
9  # 求解
10 res = optimize.linprog(-c, A, B, Aeq, Beq)
11 print(res)
12

```

```

1 con: array([3.75850107e-09])
2     fun: -201.9999999893402
3     message: 'Optimization terminated successfully.'
4     nit: 6
5     slack: array([ 93.          , 296.99999998])
6     status: 0
7     success: True
8     x: array([1.01000000e+02, 6.13324051e-10, 3.61350245e-10])
9
10 Process finished with exit code 0

```

pulp求解

- 也可以使用pulp求解，见<https://www.bilibili.com/video/BV12h411d7Dm?p=4>
- 但是稍微繁琐

pymprog求解

- 官方文档: <http://pymprog.sourceforge.net/>
<http://pymprog.sourceforge.net/intro.html#whetting>

举例

```

1 maximize 15 x + 10 y          # profit
2 S.T.
3         x          <= 3      # mountain bike limit
4         y          <= 4      # racer limit
5         x + y      <= 5      # frame limit
6         x >= 0, y >= 0       # non-negative

```

```

1 from pymprog import *
2
3 if __name__ == '__main__':
4     begin('bike production')
5     x, y = var('x, y') # variables
6     maximize(15 * x + 10 * y, 'profit')
7     x <= 3 # mountain bike limit
8     y <= 4 # racer production limit
9     x + y <= 5 # metal finishing limit
10    solve()
11
12    print('x取值: ' + str(x.primal))
13    print('y取值: ' + str(y.primal))
14    print('最优解为: ' + str(vobj()))
15

```

- res

```
1 GLPK Simplex Optimizer, v4.65
2 1 row, 2 columns, 2 non-zeros
3 * 0: obj = -0.000000000e+00 inf = 0.000e+00 (2)
4 * 2: obj = 6.500000000e+01 inf = 0.000e+00 (0)
5 OPTIMAL LP SOLUTION FOUND
6 x取值: 3.0
7 y取值: 2.0
8 最优解为: 65.0
9
10 Process finished with exit code 0
```

整数规划

cvxpy求解

- 和线性规划差不多，但是多了个约束，那就是部分变量被约束为整数
- 目前没有一种方法可以有效地求解一切整数规划。常见的整数规划求解算法有：
 - 分支定界法：可求纯或混合整数线性规划；
 - 割平面法：可求纯或混合整数线性规划；
 - 隐枚举法：用于求解0-1整数规划，有过滤隐枚举法和分支隐枚举法；
 - 匈牙利法：解决指派问题（0-1规划特殊情形）；
 - Monte Carlo法：求解各种类型规划。

举例1

$$\begin{aligned} \min \quad & z = 40x_1 + 90x_2, \\ \text{s.t.} \quad & \begin{cases} 9x_1 + 7x_2 \leq 56, \\ -7x_1 - 20x_2 \leq -70, \\ x_1, x_2 \geq 0 \text{ 为整数.} \end{cases} \end{aligned}$$

- 同理也是化成 \leq 的标准形式
- 这里的改动只需要我们输入 n, a, b, c ，以及第10行的小改动， n, a, b, c 含义和上面线性规划一样，如果有 Aeq 和 Beq 也是同理，加上即可，然后放进 $cons$ (下面第11行)里面
- 如果是求最大，第10行用 $cp.Maximize$

```
1 import cvxpy as cp
2 from numpy import array
3
4 if __name__ == '__main__':
5     n = 2 # 两个变量
6     c = array([40, 90]) # 定义目标向量
7     a = array([[9, 7], [-7, -20]]) # 定义约束矩阵
```

```

8     b = array([56, -70]) # 定义约束条件的右边向量
9     x = cp.Variable(n, integer=True) # 定义两个整数决策变量
10    obj = cp.Minimize(c * x) # 构造目标函数
11    cons = [a * x <= b, x >= 0] # 构造约束条件
12    prob = cp.Problem(obj, cons) # 构建问题模型
13    prob.solve(solver='GLPK_MI', verbose=True) # 求解问题
14    # prob.solve(solver=cp.CPLEX, verbose=True) # cp.CPLEX也可以
15    print("最优值为:", prob.value)
16    print("最优解为:", x.value)

```

- 运行结果会有警告，但是不影响结果

```

1      0: obj = 2.700000000e+02 inf = 6.250e-01 (1)
2      1: obj = 3.150000000e+02 inf = 0.000e+00 (0)
3  Long-step dual simplex will be used
4  +      1: mip = not found yet >= -inf (1; 0)
5  Solution found by heuristic: 360
6  +      2: >>>> 3.500000000e+02 >= 3.500000000e+02 0.0% (1; 0)
7  +      2: mip = 3.500000000e+02 >= tree is empty 0.0% (0; 1)
8  最优值为: 350.0
9  最优解为: [2. 3.]
10
11  Process finished with exit code 0

```

- 参考: <https://zhuanlan.zhihu.com/p/344215929>

非线性规划

- 非线性规划可分为两种，目标函数是凸函数或者是非凸函数
 - 凸函数的非线性规划：如 $f = x^2 + y^2 + x \cdot y$ ，可以使用 `scipy`
 - 非凸函数的非线性规划：如求极值，可以有如下方法
 - 纯数学方法，求导求极值
 - 神经网络，深度学习(bp算法链式求导)
 - `scipy.optimize.minimize`

```

1  fun: 求最小值的目标函数
2  args: 常数值
3  method: 求极值方法，一般默认。
4  constraints: 约束条件
5  x0: 变量的初始猜测值，注意 minimize是局部最优

```

举例

- 计算 $1/x+x$ 的最小值
 - 只需要改12行的系数，15行的初始猜测值，和8行的函数
 - 如果结果是True，则是找到局部最优解，若是False，则结果是错误的

```
1  from scipy.optimize import minimize
2  import numpy as np
3
4
5  # f = 1/x+x
6  def fun(args):
7      a = args
8      return lambda x: a / x[0] + x[0]
9
10
11 if __name__ == '__main__':
12     args = (1) # a
13     # x0 = np.asarray((1.5)) # 初始猜测值
14     # x0 = np.asarray((2.2)) # 初始猜测值
15     x0 = np.asarray((2)) # 设置初始猜测值
16
17     res = minimize(fun(args), x0, method='SLSQP')
18     print('最值:', res.fun)
19     print('是否是最优解', res.success)
20     print('取到最值时, x的值(最优解)是', res.x)
21
```

```
1  最值: 2.00000007583235
2  是否是最优解 True
3  取到最值时, x的值(最优解)是 [1.00027541]
4
5  Process finished with exit code 0
```

举例2

•

计算 $(2 + x_1)/(1 + x_2) - 3x_1 + 4x_3$ 的最小值

- x_1, x_2, x_3 的范围都在0.1到0.9 之间
- x 是变量矩阵，如 $x[0]$ 即为 x_1
- 需要变动的是函数fun，con, 27, 29, 32行， x_0 的设置要尽在要求的0.1到0.9范围内

```
1  from scipy.optimize import minimize
```

```

2 import numpy as np
3
4
5 # 计算(2+x1)/(1+x2)- 3*x1+4*x3
6 def fun(args):
7     a, b, c, d = args
8     return lambda x: (a + x[0]) / (b + x[1]) - c * x[0] + d * x[2]
9
10
11 def con(args):
12     # 约束条件 分为eq 和ineq
13     # eq表示 函数结果等于0
14     # ineq 表示 表达式大于等于0
15     x1min, x1max, x2min, x2max, x3min, x3max = args
16     cons = ({'type': 'ineq', 'fun': lambda x: x[0] - x1min},
17             {'type': 'ineq', 'fun': lambda x: -x[0] + x1max},
18             {'type': 'ineq', 'fun': lambda x: x[1] - x2min},
19             {'type': 'ineq', 'fun': lambda x: -x[1] + x2max},
20             {'type': 'ineq', 'fun': lambda x: x[2] - x3min},
21             {'type': 'ineq', 'fun': lambda x: -x[2] + x3max})
22     return cons
23
24
25 if __name__ == "__main__":
26     # 定义常量值
27     args = (2, 1, 3, 4) # a,b,c,d
28     # 设置参数范围/约束条件
29     args1 = (0.1, 0.9, 0.1, 0.9, 0.1, 0.9) # x1min, x1max, x2min,x2max
30     cons = con(args1)
31     # 设置初始猜测值
32     x0 = np.asarray((0.5, 0.5, 0.5))
33     res = minimize(fun(args), x0, method='SLSQP', constraints=cons)
34     print('最值:', res.fun)
35     print('是否是最优解', res.success)
36     print('取到最值时, x的值(最优解)是', res.x)
37

```

```

1 最值: -0.773684210526435
2 是否是最优解 True
3 取到最值时, x的值(最优解)是 [0.9 0.9 0.1]
4
5 Process finished with exit code 0

```

- 可以看出对于这类简单函数，局部最优解与真实最优解相差不大，但是对于复杂的函数，x0的初始值设置，会很大程度影响最优解的结果

数值逼近

一维和二维插值

- 参考<https://www.bilibili.com/video/BV12h411d7Dm?p=5>
- 都使得图像更加光滑

最小二乘法拟合

- 用的是

```
1 from scipy.optimize import leastsq
```

举例

- 一组数据:

```
1 X = np.array([8.19, 2.72, 6.39, 8.71, 4.7, 2.66, 3.78])
2 Y = np.array([7.01, 2.78, 6.47, 6.71, 4.1, 4.23, 4.05])
```

- 使用leastsq

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from scipy.optimize import leastsq
4 from pylab import mpl
5
6 mpl.rcParams['font.sans-serif'] = ['Microsoft YaHei'] # 指定默认字体
7 mpl.rcParams['axes.unicode_minus'] = False # 解决保存图像是负号 '-' 显示为方块的问题
8
9
10 # 计算以p为参数的直线与原始数据之间误差
11 def f(p):
12     k, b = p
13     return Y - (k * X + b)
14
15
16 if __name__ == '__main__':
17     X = np.array([8.19, 2.72, 6.39, 8.71, 4.7, 2.66, 3.78])
18     Y = np.array([7.01, 2.78, 6.47, 6.71, 4.1, 4.23, 4.05])
19     # leastsq使得f的输出数组的平方和最小, 参数初始值为[1,0]
20     r = leastsq(f, [1, 0]) # 数初始值可以随便设个合理的
21     k, b = r[0]
22     x = np.linspace(0, 10, 1000)
23     y = k * x + b
24
25     # 画散点图, s是点的大小
26     plt.scatter(X, Y, s=100, alpha=1.0, marker='o', label=u'数据点')
27     # 话拟合曲线, linewidth是线宽
28     plt.plot(x, y, color='r', linewidth=2, linestyle="--", markersize=20, label=u'拟合曲线')
```

Made By Jackpu

```

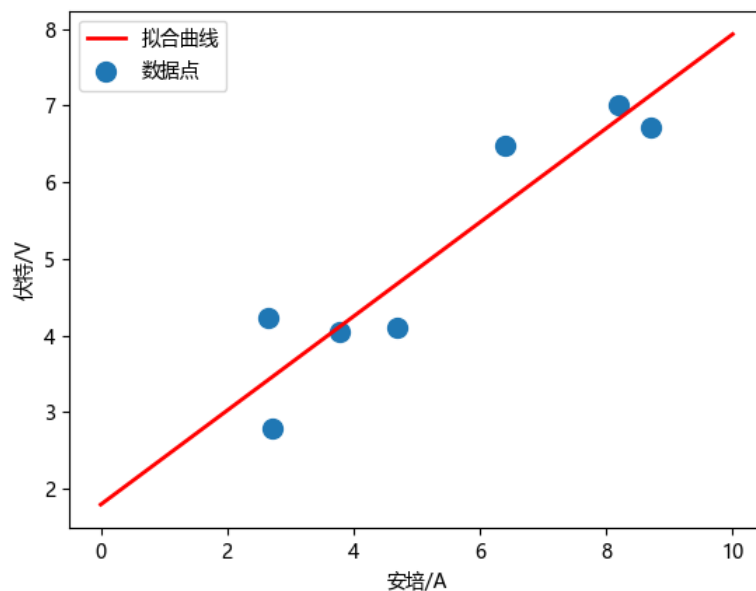
29 plt.xlabel('安培/A') # 美赛就不用中文了
30 plt.ylabel('伏特/V')
31 plt.legend(loc=0, numpoints=1) # 显示点和线的说明
32 # plt.plot(X, Y)
33 plt.show()
34
35 print('k = ', k)
36 print('b = ', b)
37

```

```

1 k = 0.6134953491930442
2 b = 1.794092543259387
3
4 Process finished with exit code 0

```



- 另外，这个线性拟合也可以使用sklearn求k和b，再去画图
 - 注意维度转换 `X = X.reshape(-1, 1)`

```

1 from sklearn import linear_model
2 import numpy as np
3
4 if __name__ == '__main__':
5     X = np.array([8.19, 2.72, 6.39, 8.71, 4.7, 2.66, 3.78])
6     Y = np.array([7.01, 2.78, 6.47, 6.71, 4.1, 4.23, 4.05])
7     X = X.reshape(-1, 1)
8     model = linear_model.LinearRegression()
9     model.fit(X, Y)
10    print('k = ', model.coef_)
11    print('b = ', model.intercept_)
12

```

```

1 k = [0.61349535]
2 b = 1.7940925542916233
3
4 Process finished with exit code 0

```

- 更多线性和非线性问题、如多元回归、逻辑回归、其他分类问题，见我之前的sklearn blog

微分方程

- 微分方程是用来描述某一类函数与其导数之间关系的方程，其解是一个符合方程的函数。微分方程按自变量个数可分为 常微分方程和偏微分方程，

前者表达通式：
$$f\left(\frac{d^n y}{dx^n}, \frac{d^{(n-1)} y}{dx^{(n-1)}}, \dots, \frac{dy}{dx}, y\right) = p(x)$$

后者表达通式：
$$\frac{\partial \phi}{\partial x} + x \frac{\partial \phi}{\partial y} = 0$$

- 建议稍微复习一下高数上册最后微分方程那章再看看会更好

解析解

- 使用sympy库，但是得到的是字符形式的格式

如下这种，如果结果是比较复杂的，可能太丑

以求解阻尼谐振子的二阶ODE为例，其表达式为：

如

$$\frac{d^2 x(t)}{dt^2} + 2\gamma\omega_0 \frac{dx(t)}{dt} + \omega_0^2 x(t) = 0$$

$$initial\ conditions. \begin{cases} x(0) = 1 \\ \left. \frac{dx(t)}{dt} \right|_{t=0} = 0 \end{cases}$$

```

1 import sympy
2
3 def apply_ics(sol, ics, x, known_params):
4     free_params = sol.free_symbols - set(known_params)
5     eqs = [(sol.lhs.diff(x, n) - sol.rhs.diff(x, n)).subs(x, 0).subs(ics) for n
6             in range(len(ics))]
7     sol_params = sympy.solve(eqs, free_params)
8     return sol.subs(sol_params)

```

```

9
10 if __name__ == '__main__':
11     sympy.init_printing() # 初始化打印环境
12     t, omega0, gamma = sympy.symbols("t, omega_0, gamma", positive=True) # 标记
    参数, 且均为正
13     x = sympy.Function('x') # 标记x是微分函数, 非变量
14     ode = x(t).diff(t, 2) + 2 * gamma * omega0 * x(t).diff(t) + omega0 ** 2 *
    x(t)
15     ode_sol = sympy.dsolve(ode) # 用diff()和dsolve得到通解
16     ics = {x(0): 1, x(t).diff(t).subs(t, 0): 0} # 将初始条件字典匹配
17     x_t_sol = apply_ics(ode_sol, ics, t, [omega0, gamma])
18     sympy.pprint(x_t_sol)
19

```

◦ 此段解释可见: <https://www.bilibili.com/video/BV12h411d7Dm?p=6>

$$\begin{aligned}
 x(t) = & \left(\frac{\gamma}{2\sqrt{\gamma^2 - 1}} + \frac{1}{2} \right) e^{\omega_0 t (-\gamma - \sqrt{\gamma^2 - 1})} + \left(\frac{\gamma}{2\sqrt{\gamma^2 - 1}} - \frac{1}{2} \right) e^{\omega_0 t (-\gamma + \sqrt{\gamma^2 - 1})} \\
 & + \left(\frac{\gamma}{2\sqrt{\gamma^2 - 1}} + \frac{1}{2} \right) e^{\omega_0 t (-\gamma - \sqrt{\gamma^2 - 1})} + \left(\frac{\gamma}{2\sqrt{\gamma^2 - 1}} - \frac{1}{2} \right) e^{\omega_0 t (-\gamma + \sqrt{\gamma^2 - 1})}
 \end{aligned}$$

Process finished with exit code 0

◦ 如果最后加上 `print(x_t_sol)`

结果为:

```

1 Eq(x(t), (-gamma/(2*sqrt(gamma**2 - 1)) + 1/2)*exp(omega_0*t*(-gamma - sqrt(gamma
    - 1)*sqrt(gamma + 1))) + (gamma/(2*sqrt(gamma**2 - 1)) + 1/2)*exp(omega_0*t*(-
    gamma + sqrt(gamma - 1)*sqrt(gamma + 1))))
2

```

• 结果较为简单的常微分方程

◦ $f(x)'' + f(x) = 0$ 二阶常系数齐次微分方程

```

1 import sympy as sy
2
3 # f(x)''+f(x)=0 二阶常系数齐次微分方程
4 def differential_equation(x, f):
5     return sy.diff(f(x), x, 2) + f(x)
6
7
8 if __name__ == '__main__':
9     x = sy.symbols('x') # 约定变量
10    f = sy.Function('f') # 约定函数
11    print(sy.dsolve(differential_equation(x, f), f(x))) # 打印
12    sy.pprint(sy.dsolve(differential_equation(x, f), f(x))) # 漂亮的打印
13

```

```

1 Eq(f(x), C1*sin(x) + C2*cos(x))
2 f(x) = C1*sin(x) + C2*cos(x)
3
4 Process finished with exit code 0

```

- 可以参考: https://blog.csdn.net/your_answer/article/details/79234275

数值解

- 当ODE(常微分方程)无法求得解析解时, 可以用scipy中的integrate.odeint求数值解来探索其解的部分性质, 并辅以可视化, 能直观地展现ODE解的函数表达

举例

- 一阶非线性常微分方程 $\frac{dy}{dx} = x - y(x)^2$
- plot_direction_field函数里面的参数意义:
 - y_x: 也就是y(x)
 - f_xy: 也就是x-y(x)^2
 - x_lim=(-5, 5), y_lim=(-5, 5)也就是在这个x, y轴的范围展示出来
- 关键需要修改的部分, 29-31行, 33行的y0需要适当调

```

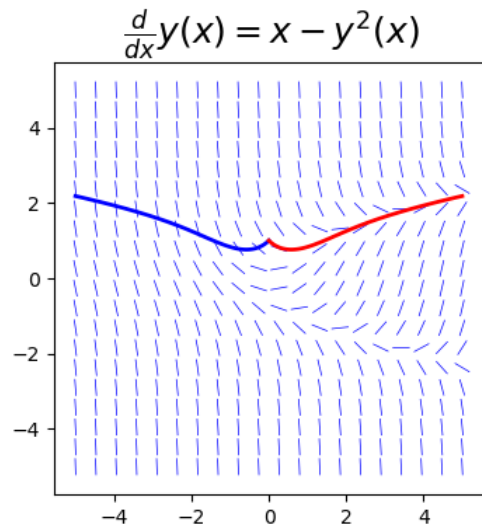
1 x = sympy.symbols('x')
2 y = sympy.Function('y')
3 f = x - y(x) ** 2

```

```

1  import numpy as np
2  from scipy import integrate
3  import matplotlib.pyplot as plt
4  import sympy
5
6
7  def plot_direction_field(x, y_x, f_xy, x_lim=(-5, 5), y_lim=(-5, 5), ax=None):
8      f_np = sympy.lambdify((x, y_x), f_xy, 'numpy')
9      x_vec = np.linspace(x_lim[0], x_lim[1], 20)
10     y_vec = np.linspace(y_lim[0], y_lim[1], 20)
11     if ax is None:
12         _, ax = plt.subplots(figsize=(4, 4))
13
14     dx = x_vec[1] - x_vec[0]
15     dy = y_vec[1] - y_vec[0]
16     for m, xx in enumerate(x_vec):
17         for n, yy in enumerate(y_vec):
18             Dy = f_np(xx, yy) * dx
19             Dx = 0.8 * dx ** 2 / np.sqrt(dx ** 2 + Dy ** 2)
20             Dy = 0.8 * Dy * dy / np.sqrt(dx ** 2 + Dy ** 2)
21             ax.plot([xx - Dx / 2, xx + Dx / 2], [yy - Dy / 2, yy + Dy / 2], 'b',
22                     lw=0.5)
23
24     ax.axis('tight')
25     ax.set_title(r'$ % s $' % (sympy.latex(sympy.Eq(y_x.diff(x), f_xy))),
26                 fontsize=18)
27     return ax
28
29 if __name__ == '__main__':
30     x = sympy.symbols('x')
31     y = sympy.Function('y')
32     f = x - y(x) ** 2
33     f_np = sympy.lambdify((y(x), x), f) # 符号表达式转隐函数
34     y0 = 1 # odeint需要给个初始值
35     xp = np.linspace(0, 5, 100)
36     yp = integrate.odeint(f_np, y0, xp) # 初始y0解f_np,x范围xp
37     xn = np.linspace(0, -5, 100)
38     yn = integrate.odeint(f_np, y0, xp)
39     fig, ax = plt.subplots(1, 1, figsize=(4, 4))
40     plot_direction_field(x, y(x), f, ax=ax) # 绘制f的场线图
41     ax.plot(xn, yn, 'b', lw=2)
42     ax.plot(xp, yp, 'r', lw=2)
43     plt.show()

```



传染病模型

- 传染病模型研究属于传染病动力学研究方向，这里只是将模型中微分方程进行了python实现
- 传染病模型包括：SI、SIS、SIR、SIRS、SEIR、SEIRS共六个模型

SI模型

- 比如病毒传染初期，没有加防疫手段，就符合SI模型
- SI模型的表达式见网上(S:易感染，I:已感染)
- 需要修改的参数

```
1 N = 10000 # N为人群总数
2 beta = 0.25 # β为传染率系数
3 gamma = 0 # gamma为恢复率系数
4 I_0 = 1 # I_0为感染者的初始人数
5 S_0 = N - I_0 # S_0为易感染者的初始人数
6 T = 150 # T为传播时间
```

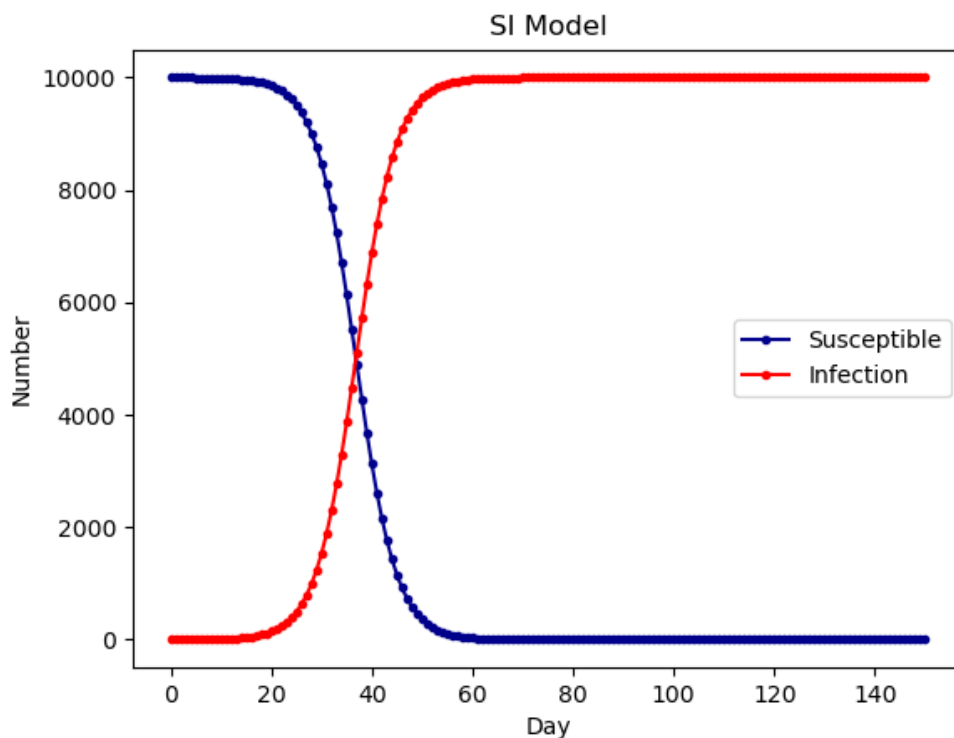
- β 为传染率系数,比如现在100个人已经传染了，在一段时间内，传染新增了25人，则 β 为0.25
- γ 为恢复率系数，一开始没有抗体都是为0的，如果不为0，比如是开始有100人感染，在一个传播时间T后，治愈了6个人，则 γ 取0.06
- I_0 为感染者的初始人数
- S_0 为易感染者的初始人数，这个要看情况，如果都不加干预，那就是 $N - I_0$ ，一般看情况需要再考虑其他因素(交通，社交群体，航线等)， S_0 考虑的越多，则越完备
- Susceptible易感染的，Infection已经感染的
- code

```
1 import numpy as np
2 import scipy.integrate as spi
3 import matplotlib.pyplot as plt
4
5 N = 10000 # N为人群总数
```

```

6  beta = 0.25 #  $\beta$ 为传染率系数
7  gamma = 0 # gamma为恢复率系数
8  I_0 = 1 # I_0为感染者的初始人数
9  S_0 = N - I_0 # S_0为易感染者的初始人数
10 T = 150 # T为传播时间
11 INI = (S_0, I_0) # INI为初始状态下的数组
12
13
14 def funcSI(inivalue, _):
15     Y = np.zeros(2)
16     X = inivalue
17     Y[0] = -(beta * X[0] * X[1]) / N + gamma * X[1] # 易感个体变化
18     Y[1] = (beta * X[0] * X[1]) / N - gamma * X[1] # 感染个体变化
19     return Y
20
21
22 if __name__ == '__main__':
23     T_range = np.arange(0, T + 1)
24     RES = spi.odeint(funcSI, INI, T_range)
25     plt.plot(RES[:, 0], color='darkblue', label='Susceptible', marker='.')
26     plt.plot(RES[:, 1], color='red', label='Infection', marker='.')
27     plt.title('SI Model')
28     plt.legend()
29     plt.xlabel('Day')
30     plt.ylabel('Number')
31     plt.show()
32

```

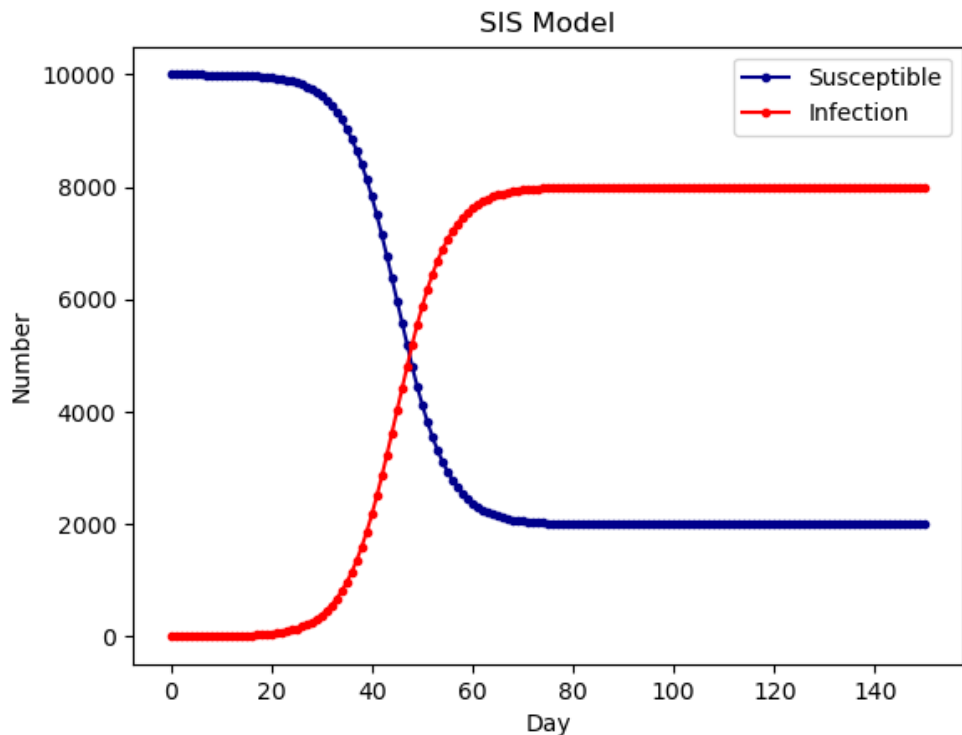


- 可以看到10000个人大概在60天左右就全部感染了

SIS模型

- 与SI区别不大，区别在于7行的gamma有初始值，以及17行的公式改变了

```
1 import numpy as np
2 import scipy.integrate as spi
3 import matplotlib.pyplot as plt
4
5 N = 10000 # N为人群总数
6 beta = 0.25 #  $\beta$ 为传染率系数
7 gamma = 0.05 # gamma为恢复率系数
8 I_0 = 1 # I_0为感染者的初始人数
9 S_0 = N - I_0 # S_0为易感染者的初始人数
10 T = 150 # T为传播时间
11 INI = (S_0, I_0) # INI为初始状态下的数组
12
13
14 def funcSI(inivalue, _):
15     Y = np.zeros(2)
16     X = inivalue
17     Y[0] = -(beta * X[0]) / N * X[1] + gamma * X[1] # 易感个体变化
18     Y[1] = (beta * X[0] * X[1]) / N - gamma * X[1] # 感染个体变化
19     return Y
20
21
22 if __name__ == '__main__':
23     T_range = np.arange(0, T + 1)
24     RES = spi.odeint(funcSI, INI, T_range)
25     plt.plot(RES[:, 0], color='darkblue', label='Susceptible', marker='.')
26     plt.plot(RES[:, 1], color='red', label='Infection', marker='.')
27     plt.title('SIS Model')
28     plt.legend()
29     plt.xlabel('Day')
30     plt.ylabel('Number')
31     plt.show()
32
```



- 可以看到60-80天之间，逐渐稳定，有的人治愈后(获得抗体)活了下来，有的没治愈的就死了

SIR模型

- 多了 R_0 为治愈者的初始人数，即刚开始注射疫苗恢复的人
- 表达式也改变
- 注意恢复治愈包括自身产生抗体以及通过医疗手段获得抗体两种

```

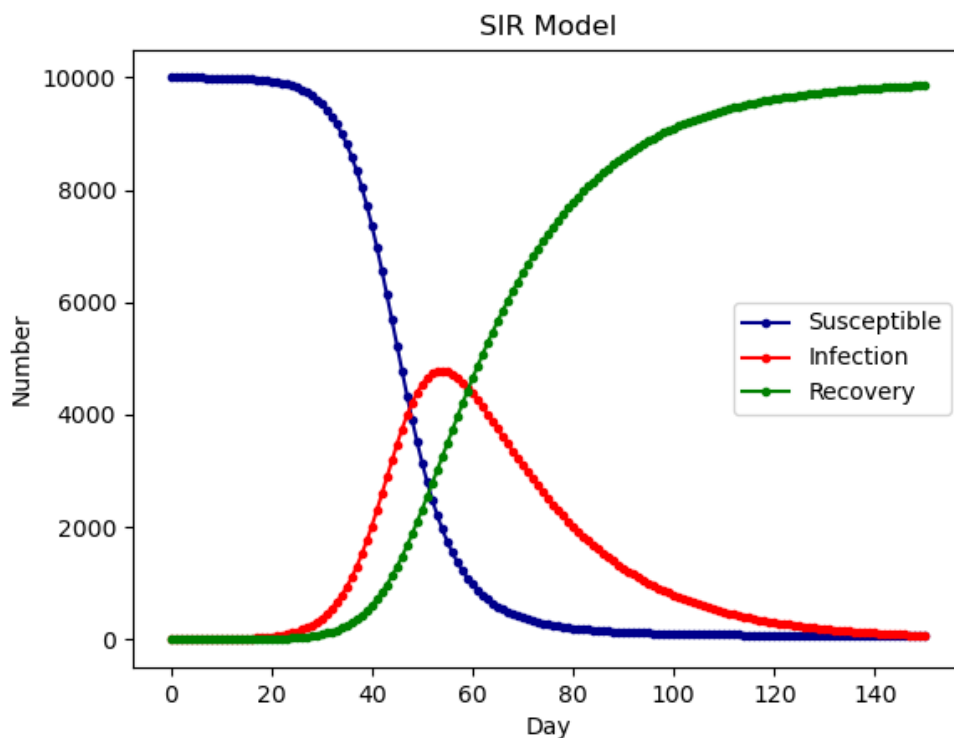
1  import numpy as np
2  import scipy.integrate as spi
3  import matplotlib.pyplot as plt
4
5  N = 10000 # N为人群总数
6  beta = 0.25 #  $\beta$ 为传染率系数
7  gamma = 0.05 # gamma为恢复率系数
8  I_0 = 1 # I_0为感染者的初始人数
9  R_0 = 0 # R_0为治愈者的初始人数
10 S_0 = N - I_0 - R_0 # S_0为易感染者的初始人数
11 T = 150 # T为传播时间
12 INI = (S_0, I_0, R_0) # INI为初始状态下的数组
13
14
15 def funcSIR(invalue, _):
16     Y = np.zeros(3)
17     X = invalue
18     Y[0] = -(beta * X[0] * X[1]) / N # 易感个体变化
19     Y[1] = (beta * X[0] * X[1]) / N - gamma * X[1] # 感染个体变化

```

```

20     Y[2] = gamma * X[1] # 治愈个体变化
21     return Y
22
23
24 if __name__ == '__main__':
25     T_range = np.arange(0, T + 1)
26     RES = spi.odeint(funcSIR, INI, T_range)
27     plt.plot(RES[:, 0], color='darkblue', label='Susceptible', marker='.')
28     plt.plot(RES[:, 1], color='red', label='Infection', marker='.')
29     plt.plot(RES[:, 2], color='green', label='Recovery', marker='.')
30     plt.title('SIR Model')
31     plt.legend()
32     plt.xlabel('Day')
33     plt.ylabel('Number')
34     plt.show()
35

```



- 可以看到感染人数出现峰值，在前期已经开始使用抗体，病人逐渐治愈，最后所有人都恢复健康

SIRS模型

- 多了Ts为抗体持续时间，也就是说有了抗体一段时间后，抗体失效，又变成了易感染人群
- 公式也改变

```

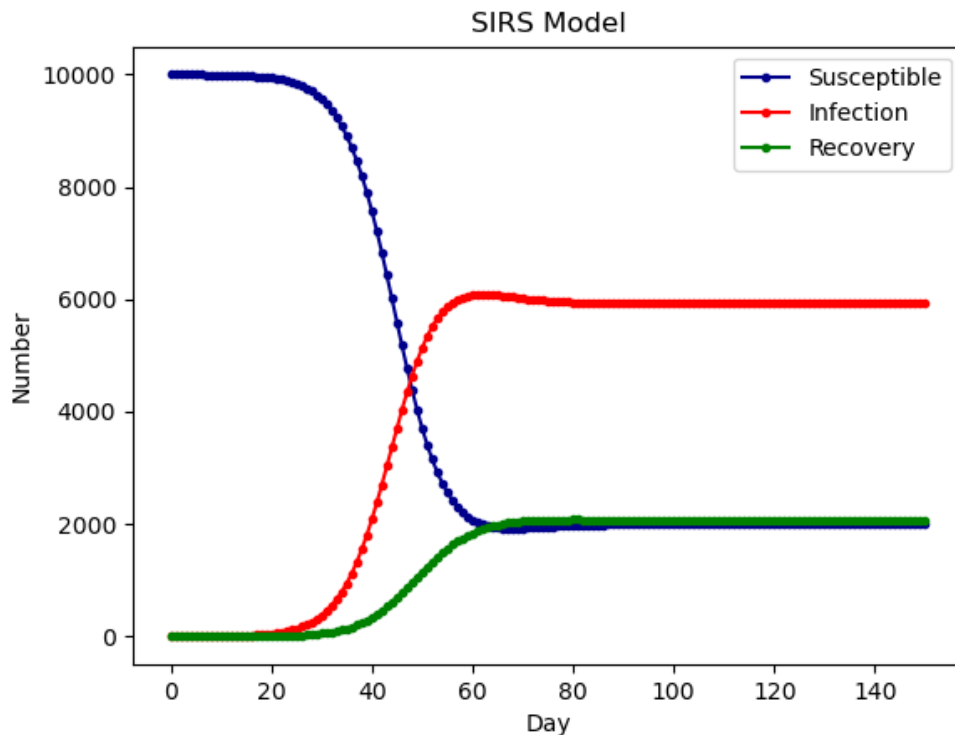
1 import numpy as np
2 import scipy.integrate as spi
3 import matplotlib.pyplot as plt

```

```

4
5 N = 10000 # N为人群总数
6 beta = 0.25 #  $\beta$ 为传染率系数
7 gamma = 0.05 # gamma为恢复率系数
8 Ts = 7 # Ts为抗体持续时间
9 I_0 = 1 # I_0为感染者的初始人数
10 R_0 = 0 # R_0为治愈者的初始人数
11 S_0 = N - I_0 - R_0 # S_0为易感染者的初始人数
12 T = 150 # T为传播时间
13 INI = (S_0, I_0, R_0) # INI为初始状态下的数组
14
15
16 def funcSIRS(inivalue, _):
17     Y = np.zeros(3)
18     X = inivalue
19     Y[0] = -(beta * X[0] * X[1]) / N + X[2] / Ts # 易感个体变化
20     Y[1] = (beta * X[0] * X[1]) / N - gamma * X[1] # 感染个体变化
21     Y[2] = gamma * X[1] - X[2] / Ts # 治愈个体变化
22     return Y
23
24
25 if __name__ == '__main__':
26     T_range = np.arange(0, T + 1)
27     RES = spi.odeint(funcSIRS, INI, T_range)
28     plt.plot(RES[:, 0], color='darkblue', label='Susceptible', marker='.')
29     plt.plot(RES[:, 1], color='red', label='Infection', marker='.')
30     plt.plot(RES[:, 2], color='green', label='Recovery', marker='.')
31     plt.title('SIRS Model')
32     plt.legend()
33     plt.xlabel('Day')
34     plt.ylabel('Number')
35     plt.show()
36

```



- 最终达到一个平衡

SEIR模型

- 考虑了病毒的潜伏期(潜伏人群E)，也就是感染病毒后，过了潜伏期就是感染人群了
- 多了E_0为潜伏者的初始人数，如果是0，那说明开始时有人感染，但是还没有发病，此时这类人不是易感染，但是他们携带病毒
- 只有经过潜伏期，才能被传染
- 公式有所变化

```

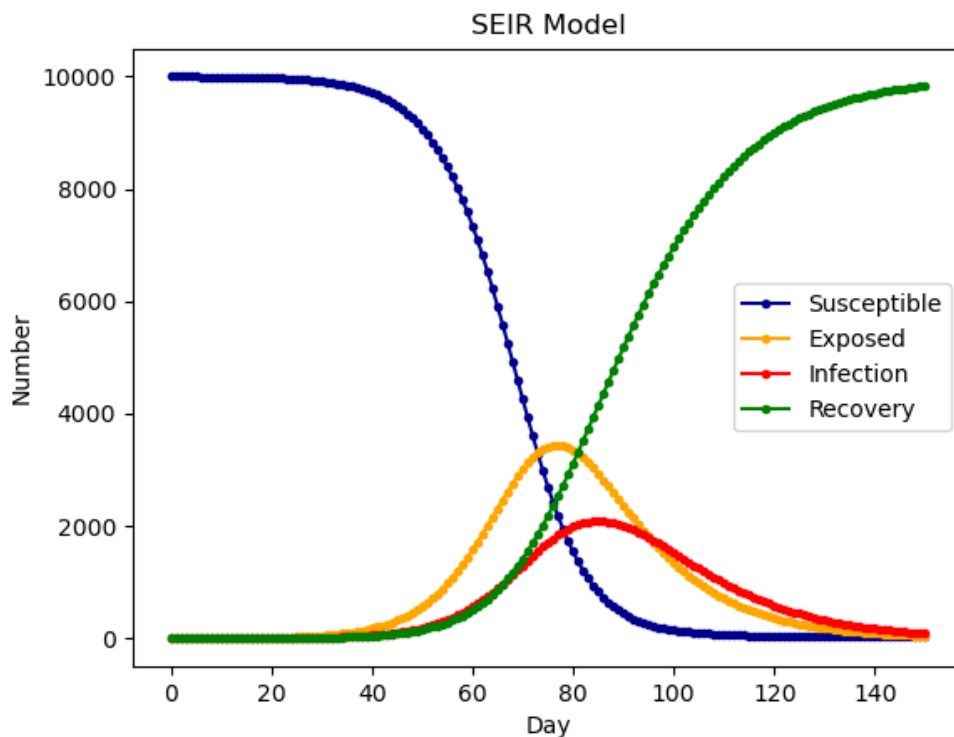
1  import numpy as np
2  import scipy.integrate as spi
3  import matplotlib.pyplot as plt
4
5  N = 10000 # N为人群总数
6  beta = 0.6 # β为传染率系数
7  gamma = 0.1 # gamma为恢复率系数
8  Te = 14 # Te为疾病潜伏期
9  I_0 = 1 # I_0为感染者的初始人数
10 E_0 = 0 # E_0为潜伏者的初始人数
11 R_0 = 0 # R_0为治愈者的初始人数
12 S_0 = N - I_0 - R_0 - E_0 # S_0为易感染者的初始人数
13 T = 150 # T为传播时间
14 INI = (S_0, E_0, I_0, R_0) # INI为初始状态下的数组
15
16
17 def funcSEIR(inivalue, _):

```

```

18 Y = np.zeros(4)
19 X = inivalue
20 Y[0] = -(beta * X[0] * X[2]) / N # 易感个体变化
21 Y[1] = (beta * X[0] * X[2] / N - X[1] / Te) # 潜伏个体变化
22 Y[2] = X[1] / Te - gamma * X[2] # 感染个体变化
23 Y[3] = gamma * X[2] # 治愈个体变化
24 return Y
25
26
27 if __name__ == '__main__':
28     T_range = np.arange(0, T + 1)
29     RES = spi.odeint(funcSEIR, INI, T_range)
30     plt.plot(RES[:, 0], color='darkblue', label='Susceptible', marker=
31             '.')
32     plt.plot(RES[:, 1], color='orange', label='Exposed', marker='.')
33     plt.plot(RES[:, 2], color='red', label='Infection', marker='.')
34     plt.plot(RES[:, 3], color='green', label='Recovery', marker='.')
35     plt.title('SEIR Model')
36     plt.legend()
37     plt.xlabel('Day')
38     plt.ylabel('Number')
39     plt.show()
40

```

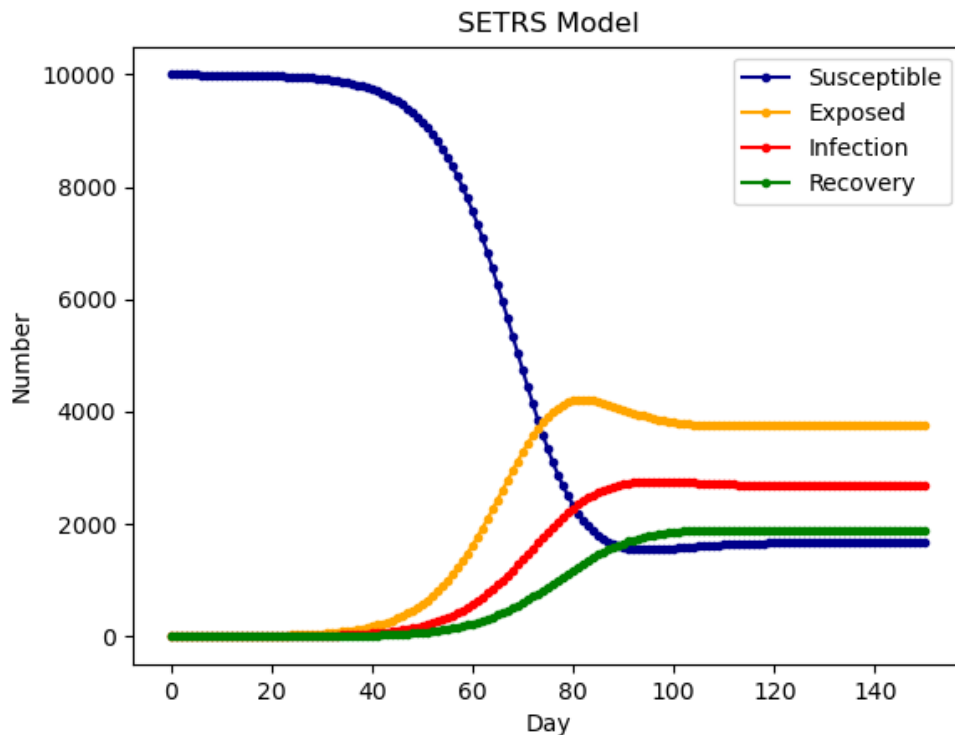


- 潜伏期的人会比感染人群先到达峰值，最终都可以治愈

SEIRS模型

- 考虑了抗体持续时间
- 一般多了潜伏期的话，传染率系数会有所增加，上面的SEIR也是同理

```
1 import numpy as np
2 import scipy.integrate as spi
3 import matplotlib.pyplot as plt
4
5 N = 10000 # N为人群总数
6 beta = 0.6 #  $\beta$ 为传染率系数
7 gamma = 0.1 # gamma为恢复率系数
8 Ts = 7 # Ts为抗体持续时间
9 Te = 14 # Te为疾病潜伏期
10 I_0 = 1 # I_0为感染者的初始人数
11 E_0 = 0 # E_0为潜伏者的初始人数
12 R_0 = 0 # R_0为治愈者的初始人数
13 S_0 = N - I_0 - R_0 - E_0 # S_0为易感染者的初始人数
14 T = 150 # T为传播时间
15 INI = (S_0, E_0, I_0, R_0) # INI为初始状态下的数组
16
17
18 def funcSEIRS(invalue, _):
19     Y = np.zeros(4)
20     X = invalue
21     Y[0] = -(beta * X[0] * X[2]) / N + X[3] / Ts # 易感个体变化
22     Y[1] = (beta * X[0] * X[2] / N - X[1] / Te) # 潜伏个体变化
23     Y[2] = X[1] / Te - gamma * X[2] # 感染个体变化
24     Y[3] = gamma * X[2] - X[3] / Ts # 治愈个体变化
25     return Y
26
27
28 if __name__ == '__main__':
29     T_range = np.arange(0, T + 1)
30     RES = spi.odeint(funcSEIRS, INI, T_range)
31     plt.plot(RES[:, 0], color='darkblue', label='Susceptible', marker='.')
32     plt.plot(RES[:, 1], color='orange', label='Exposed', marker='.')
33     plt.plot(RES[:, 2], color='red', label='Infection', marker='.')
34     plt.plot(RES[:, 3], color='green', label='Recovery', marker='.')
35     plt.title('SEIRS Model')
36     plt.legend()
37     plt.xlabel('Day')
38     plt.ylabel('Number')
39     plt.show()
40
```



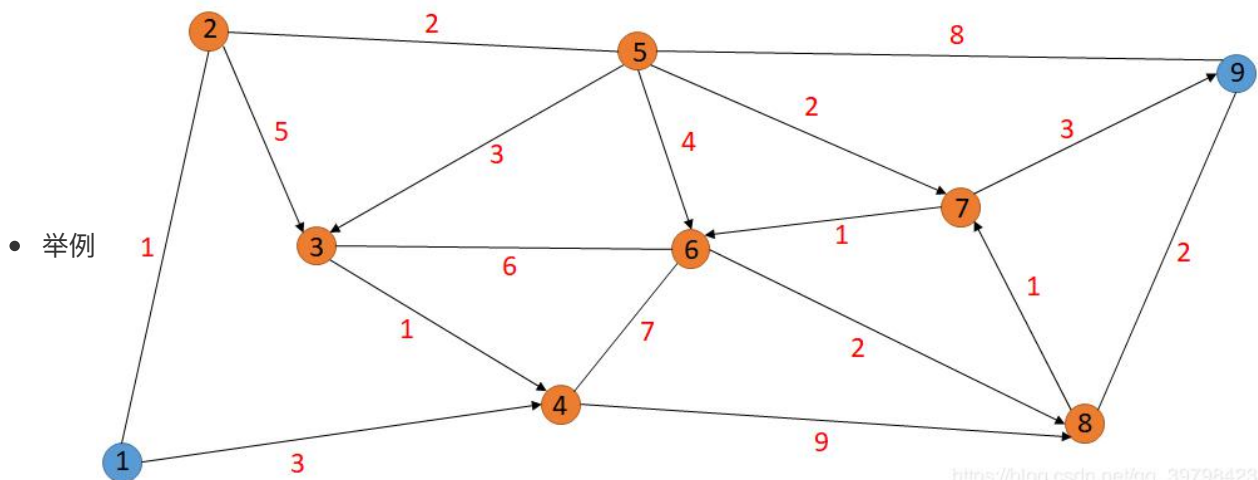
- 潜伏期的人先到峰值，然后是易感染者，然后是治愈者，他们最终会达到平衡稳定

图论

Dijkstra

解法1(常用)

- 只需要给出带有权值的邻接矩阵即可求出最短路径和最短距离
- 需要改变的第54行邻接矩阵的权值和65行的起点和终点，注意21行是从0还是1开始
- 有向边和无向边的混合均可使用
- `g = defaultdict(list)` 是得到一个元素全是list类型的字典



- 举例

https://blog.csdn.net/qz_39798423

Made By Jackp


```

1  # dijkstra
2  from collections import defaultdict
3  from heapq import *
4
5  inf = 99999 # 不连通值
6
7
8  def init_graph(mtx_graph):
9      m_n = len(mtx_graph) # 带权连接矩阵的阶数
10     edges = [] # 保存连通的两个点之间的距离(点A、点B、距离)
11     for i in range(m_n):
12         for j in range(m_n):
13             if i != j and mtx_graph[i][j] != inf:
14                 edges.append((i, j, mtx_graph[i][j]))
15
16     return edges
17
18
19  def dijkstra(edges, from_node, to_node):
20     go_path = []
21     to_node = to_node - 1 # 看情况, 如果是从1开始的就减一
22     g = defaultdict(list)
23     for l, r, c in edges:
24         g[l].append((c, r))
25     q, seen = [(0, from_node - 1, ())], set()
26     while q:
27         (cost, v1, path) = heappop(q) # 堆弹出当前路径最小成本
28         if v1 not in seen:
29             seen.add(v1)
30             path = (v1, path)
31             if v1 == to_node:
32                 break
33             for c, v2 in g.get(v1, ()):
34                 if v2 not in seen:
35                     heappush(q, (cost + c, v2, path))
36     if v1 != to_node: # 无法到达
37         return float('inf'), []
38
39     if len(path) > 0:
40         left = path[0]
41         go_path.append(left)
42         right = path[1]
43         while len(right) > 0:
44             left = right[0]
45             go_path.append(left)
46             right = right[1]
47         go_path.reverse() # 逆序变换
48         for i in range(len(go_path)): # 标号加1
49             go_path[i] = go_path[i] + 1
50     return cost, go_path

```

```

51
52
53 if __name__ == '__main__':
54     mtx_graph = [[0, 1, inf, 3, inf, inf, inf, inf, inf],
55                  [1, 0, 5, inf, 2, inf, inf, inf, inf],
56                  [inf, inf, 0, 1, inf, 6, inf, inf, inf],
57                  [inf, inf, inf, 0, inf, 7, inf, 9, inf],
58                  [inf, 2, 3, inf, 0, 4, 2, inf, 8],
59                  [inf, inf, 6, 7, inf, 0, inf, 2, inf],
60                  [inf, inf, inf, inf, inf, 1, 0, inf, 3],
61                  [inf, inf, inf, inf, inf, inf, 1, 0, 2],
62                  [inf, inf, inf, inf, 8, inf, 3, 2, 0]]
63
64     edges = init_graph(mtx_graph)
65     length, path = dijkstra(edges, 1, 9)
66     print('最短距离为: ' + str(length))
67     print('前进路径为: ' + str(path))
68

```

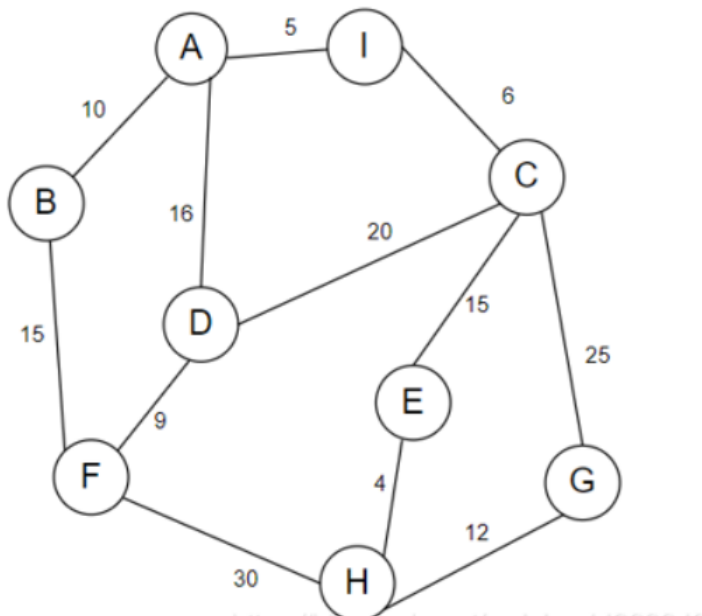
```

1  最短距离为: 8
2  前进路径为: [1, 2, 5, 7, 9]
3
4  Process finished with exit code 0

```

解法2

- 输入是一个包含每个点与其他点联系和权值的字典
- 下面是无向边的例子，适合无向边
- 需要修改的是51的联系，62的起点，66行的起点和终点



```

1  import heapq
2  Max = 99999999
3

```

```

4
5 # 有点像BFS的思想
6 def dijktestra(graph, start):
7     queue = [] # 优先队列
8     heapq.heappush(queue, (0, start))
9     visited = set()
10    path = {start: None} # 记录该点的上一个点
11
12    # 先把一开始到达的所有路径距离设最大
13    distance = {start: 0}
14    for vertex in graph:
15        if vertex != start:
16            distance[vertex] = Max
17
18    while len(queue):
19        # 取出的当前在queue的第一个点
20        pair = heapq.heappop(queue)
21        dist = pair[0]
22        vertex = pair[1]
23        visited.add(vertex)
24
25        # 该点的所有连接点
26        nodes = graph[vertex].keys()
27        for v in nodes:
28            if v not in visited and dist + graph[vertex][v] < distance[v]:
29                heapq.heappush(queue, (dist + graph[vertex][v], v)) # 优先队列会自动把
值最小的放在前面
30                path[v] = vertex # 记录上一个点
31                distance[v] = dist + graph[vertex][v] # 更新最小值
32
33    return path, distance
34
35
36 def show_path(path, start, end):
37     shortest_path = []
38     vertex = end
39     while vertex != path[start]:
40         vertex = path[vertex]
41         shortest_path.append(vertex)
42
43     shortest_path.reverse()
44     shortest_path.pop(0)
45     shortest_path.append(end)
46
47     return shortest_path
48
49
50 if __name__ == '__main__':
51     graph = {
52         'A': {'B': 10, 'D': 16, 'I': 5},

```

```

53         'B': {'A': 10, 'F': 15},
54         'C': {'D': 20, 'E': 15, 'I': 6},
55         'D': {'A': 16, 'C': 20, 'F': 9},
56         'E': {'C': 15, 'H': 4},
57         'F': {'B': 15, 'H': 30},
58         'G': {'C': 25, 'H': 12},
59         'H': {'E': 4, 'F': 9, 'G': 12},
60         'I': {'A': 5, 'C': 6}
61     }
62     path, distance = dijktestra(graph, 'A')
63     print(path)
64     print(distance)
65
66     shortest_path = show_path(path, 'A', 'H')
67     print('shortest_path:', shortest_path)
68

```

```

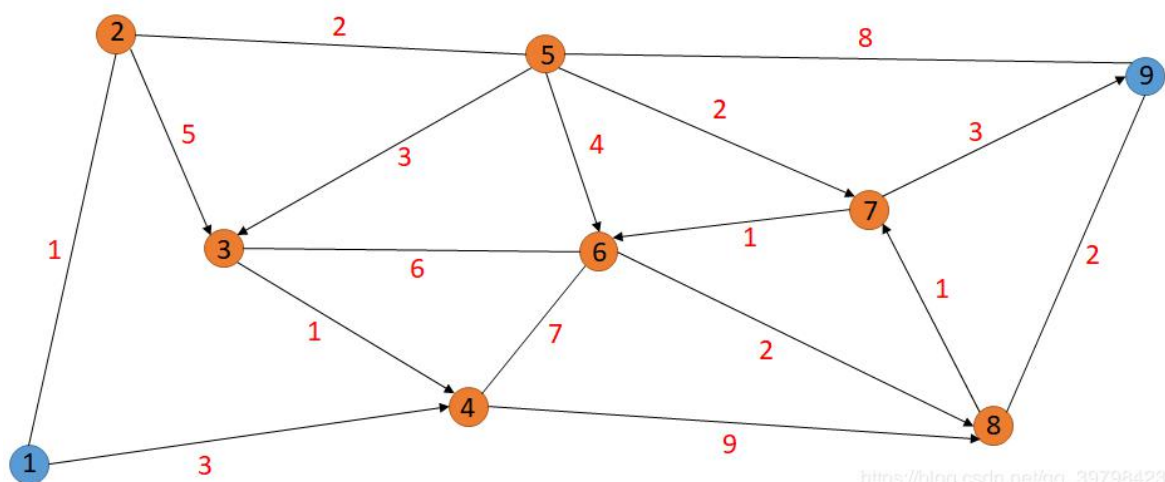
1  该点的上一个点: {'A': None, 'B': 'A', 'D': 'A', 'I': 'A', 'C': 'I', 'F': 'B', 'E': 'C',
2  起点到其他各个点的最小距离: {'A': 0, 'B': 10, 'C': 11, 'D': 16, 'E': 26, 'F': 25, 'G':
3  shortest_path: ['A', 'I', 'C', 'E', 'H']

```

Floyd

- 通过动态规划求解多源最短路径问题

举例



https://blog.csdn.net/qq_39798423

- 图和上面的图一样，求出从每一个点到其他点的最短距离和路径

```

1  import numpy as np
2
3  inf = 99999 # 不连通值
4
5
6  def floyd(graph):
7      N = len(graph)
8      A = np.array(graph)
9      path = np.zeros((N, N))
10     for i in range(0, N):
11         for j in range(0, N):
12             if A[i][j] != inf:
13                 path[i][j] = j
14
15     for k in range(0, N):
16         for i in range(0, N):
17             for j in range(0, N):
18                 if A[i][k] + A[k][j] < A[i][j]:
19                     A[i][j] = A[i][k] + A[k][j]
20                     path[i][j] = path[i][k]
21
22     for i in range(0, N):
23         for j in range(0, N):
24             path[i][j] = path[i][j] + 1
25
26     print('距离 = \n', A)
27     print('路径 = \n', path)
28
29
30 if __name__ == '__main__':
31     mtx_graph = [[0, 1, inf, 3, inf, inf, inf, inf, inf],
32                  [1, 0, 5, inf, 2, inf, inf, inf, inf],
33                  [inf, inf, 0, 1, inf, 6, inf, inf, inf],
34                  [inf, inf, inf, 0, inf, 7, inf, 9, inf],
35                  [inf, 2, 3, inf, 0, 4, 2, inf, 8],
36                  [inf, inf, 6, 7, inf, 0, inf, 2, inf],
37                  [inf, inf, inf, inf, inf, inf, 1, 0, inf, 3],
38                  [inf, inf, inf, inf, inf, inf, inf, 1, 0, 2],
39                  [inf, inf, inf, inf, 8, inf, 3, 2, 0]]
40     floyd(mtx_graph)
41

```

```

1  距离 =
2  [[ 0  1  6  3  3  6  5  8  8]
3   [ 1  0  5  4  2  5  4  7  7]
4   [21 20  0  1 18  6  9  8 10]
5   [22 21 13  0 19  7 10  9 11]
6   [ 3  2  3  4  0  3  2  5  5]
7   [15 14  6  7 12  0  3  2  4]

```

```

8  [14 13 7 8 11 1 0 3 3]
9  [13 12 8 9 10 2 1 0 2]
10 [11 10 10 11 8 4 3 2 0]]
11 路径 =
12 [[1. 2. 2. 4. 2. 2. 2. 2. 2.]
13 [1. 2. 3. 1. 5. 5. 5. 5. 5.]
14 [6. 6. 3. 4. 6. 6. 6. 6. 6.]
15 [8. 8. 6. 4. 8. 6. 8. 8. 8.]
16 [2. 2. 3. 3. 5. 7. 7. 7. 7.]
17 [8. 8. 3. 4. 8. 6. 8. 8. 8.]
18 [9. 9. 6. 6. 9. 6. 7. 6. 9.]
19 [9. 9. 7. 7. 9. 7. 7. 8. 9.]
20 [5. 5. 7. 7. 5. 7. 7. 8. 9.]]
21
22 Process finished with exit code 0

```

- 距离，比如1到9的距离为8，即从第1行看到第9列
- 怎么看最短路径看解释见下图

```

1 距离 =
2  [[ 0  1  6  3  3  6  5  8  8]
3  [ 1  0  5  4  2  5  4  7  7]
4  [21 20  0  1 18  6  9  8 10]
5  [22 21 13  0 19  7 10  9 11]
6  [ 3  2  3  4  0  3  2  5  5]
7  [15 14  6  7 12  0  3  2  4]
8  [14 13  7  8 11  1  0  3  3]
9  [13 12  8  9 10  2  1  0  2]
10 [11 10 10 11  8  4  3  2  0]]
11 路径 =
12 [1. 2. 2. 4. 2. 2. 2. 2. 2.]
13 [1. 2. 3. 1. 5. 5. 5. 5. 5.]
14 [6. 6. 3. 4. 6. 6. 6. 6. 6.]
15 [8. 8. 6. 4. 8. 6. 8. 8. 8.]
16 [2. 2. 3. 3. 5. 7. 7. 7. 7.]
17 [8. 8. 3. 4. 8. 6. 8. 8. 8.]
18 [9. 9. 6. 6. 9. 6. 7. 6. 9.]
19 [9. 9. 7. 7. 9. 7. 7. 8. 9.]
20 [5. 5. 7. 7. 5. 7. 7. 8. 9.]
21
22 Process finished with exit code 0

```

1 2 5 7 9

机场航线设计

- 图的可视化

- 数据清洗分析，可参考Kaggle练习
- 城市可作为图节点
- 这种一般考虑
 - 找到最密集的点，作为交通枢纽，考虑其他成本、时效性、盈利因素之类的...

回归

- 多元回归、逻辑回归见我之前的blog: <https://www.cnblogs.com/jmchen/p/13550573.html>

差分方程

递推关系

- 差分方程建模的关键在于如何得到第n组数据与第n+1组数据之间的关系

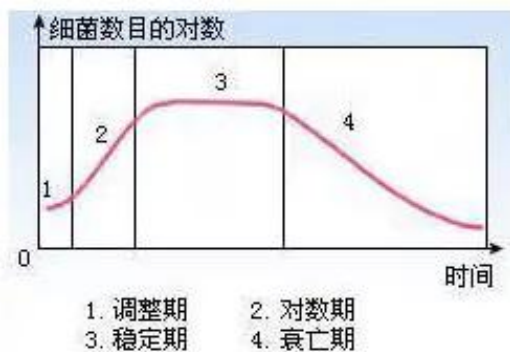
举例

酵母菌生长模型

- 相类比的还有比如兔子(其他生物)繁殖模型等
- 如图所示我们用培养基培养细菌时，其数量变化通常会经历这四个时期。

这个模型针对前三个时期建一个大致的模型：

调整期、对数期、稳定期



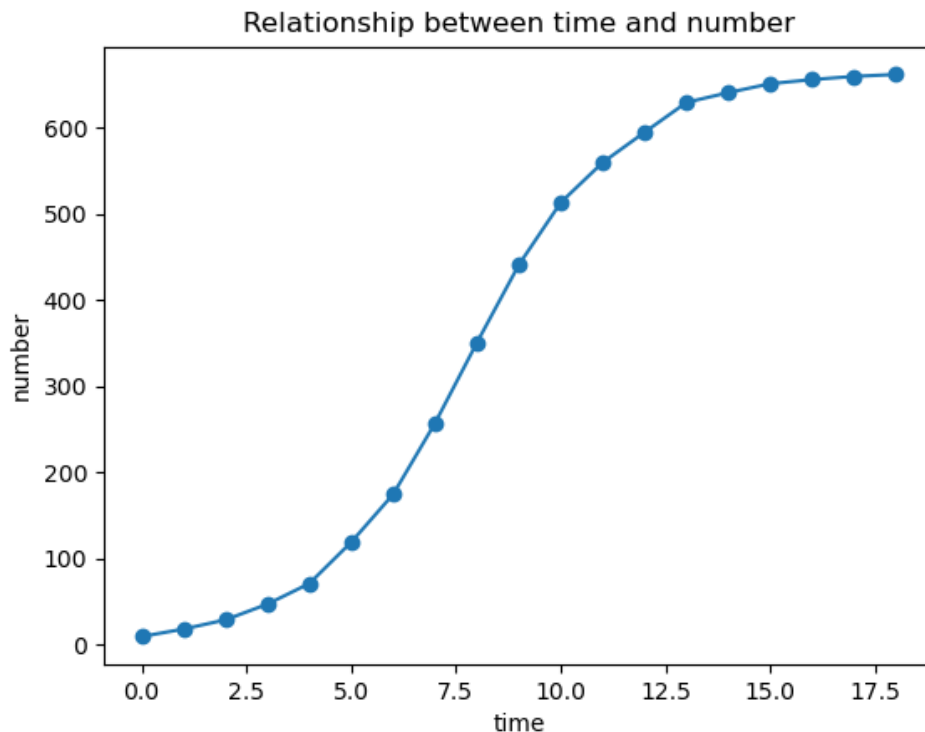
- 数据可以也从文件读入，这里直接写了

```
1 import matplotlib.pyplot as plt
2
3 if __name__ == '__main__':
4     time = [i for i in range(0, 19)]
5     number = [9.6, 18.3, 29, 47.2, 71.1, 119.1, 174.6,
6               257.3, 350.7, 441.0, 513.3, 559.7, 594.8,
7               629.4, 640.8, 651.1, 655.9, 659.6, 661.8]
8     plt.title('Relationship between time and number') # 创建标题
9     plt.xlabel('time') # X轴标签
10    plt.ylabel('number') # Y轴标签
```

```

11 plt.scatter(time, number)
12 plt.plot(time, number) # 画图
13 plt.show() # 显示
14

```



➤ 分析:

酵母菌数量增长有一个这样的规律：当某些资源只能支撑某个最大限度的种群数量，而不能支持种群数量的无限增长，当接近这个最大值时，种群数量的增长速度就会慢下来。

- 1. 两个观测点的值差 Δp 来表征增长速度
- 2. Δp 与目前的种群数量有关，数量越大，增长速度越快
- 3. Δp 还与剩余的未分配的资源量有关，资源越多，增长速度越快
- 4. 然后以极限总群数量与现有种群数量的差值表征剩余资源量

➤ 模型:

$$\Delta p = p_{n+1} - p_n = k(665 - p_n)p_n$$

- Δp : 因为横坐标间隔是1，所以相邻纵坐标之差可以当成增速
- 665是极限总群数量
- 要求的是k, 然后预测下一年
- 需要修改的是4, 17, 18行，如果不只是算下一年，那要改40, 46行

```

1 import matplotlib.pyplot as plt
2 import numpy as np
3
4 Max = 665
5

```



```

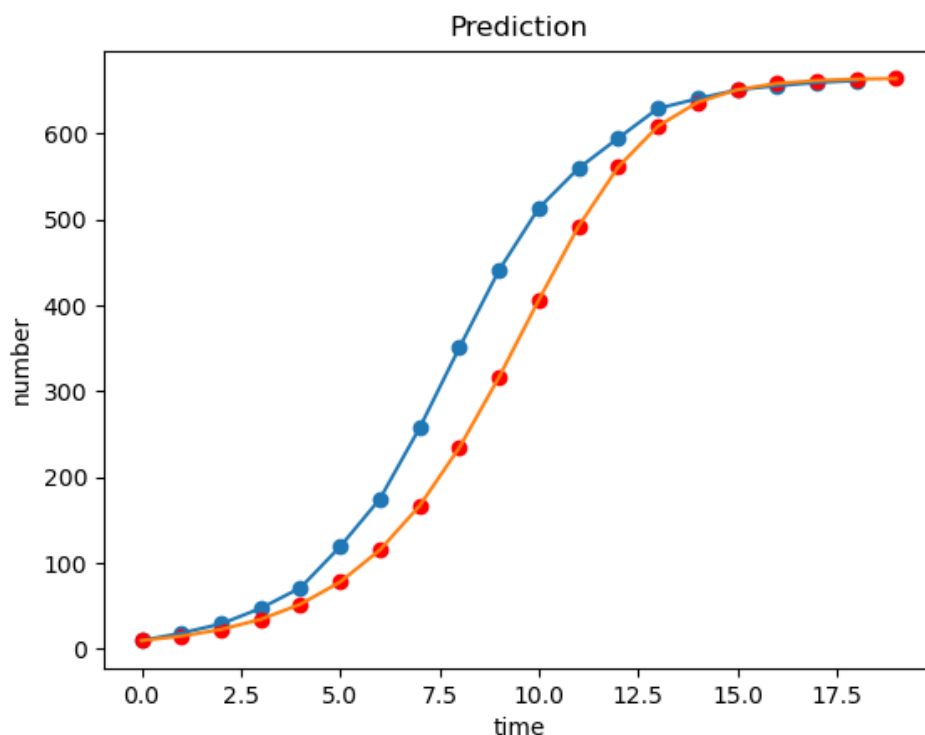
6
7 # 获取相邻纵坐标的差值
8 def get_delta(y_num: list):
9     delta_y = []
10    for i in range(len(y_num) - 1):
11        delta_y.append(y_num[i + 1] - y_num[i])
12
13    return delta_y
14
15
16 if __name__ == '__main__':
17     time = [_ for _ in range(0, 19)]
18     number = [9.6, 18.3, 29, 47.2, 71.1, 119.1, 174.6,
19               257.3, 350.7, 441.0, 513.3, 559.7, 594.8,
20               629.4, 640.8, 651.1, 655.9, 659.6, 661.8]
21
22     plt.title('Relationship between time and number') # 创建标题
23     plt.xlabel('time') # X轴标签
24     plt.ylabel('number') # Y轴标签
25     plt.scatter(time, number)
26     plt.plot(time, number) # 画图
27     # plt.show() # 显示, 注释掉后, 实际曲线和预测曲线泛在同一个图里面对比
28
29     delta_p = get_delta(number)
30     number.pop(-1)
31     pn = np.array(number)
32     f = pn * (Max - pn)
33     res = np.polyfit(f, delta_p, 1)
34     print(res)
35     print('k = ', res[0])
36
37     # 预测
38     p0 = number[0]
39     p_list = []
40     for i in range(len(time) + 1):
41         p_list.append(p0)
42         p0 = res[0] * (Max - p0) * p0 + p0
43     plt.xlabel('time') # X轴标签
44     plt.ylabel('number') # Y轴标签
45     plt.title('Prediction') # 创建标题
46     plt.scatter([_ for _ in range(0, len(time) + 1)], p_list, c='r')
47     plt.plot(p_list)
48     plt.show()
49

```

```

1 [ 0.00081448 -0.30791574]
2 k = 0.0008144797937893836
3
4 Process finished with exit code 0

```



显式差分

- 热传导方程，见<https://www.bilibili.com/video/BV12h411d7Dm?p=8>

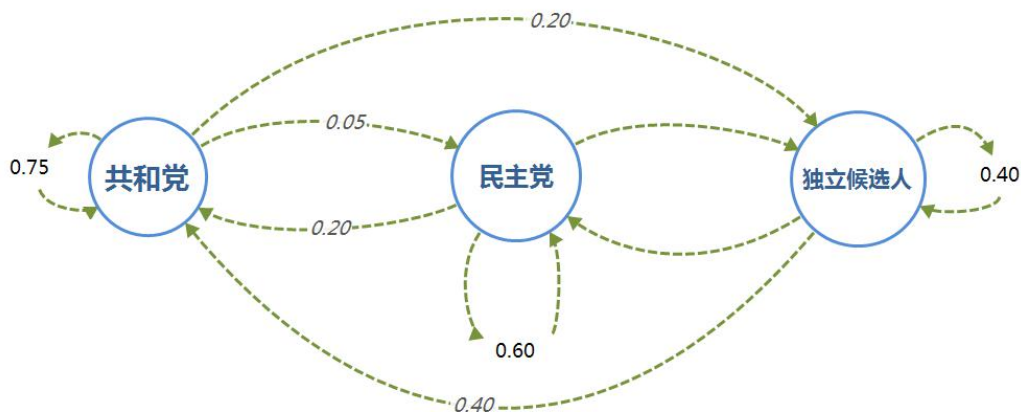
马尔科夫链

选举投票预测

- 马尔科夫链是由具有以下性质的一系列事件构成的过程：
 - 一个事件有有限多个结果，称为状态，该过程总是这些状态中的一个；
 - 在过程的每个阶段或者时段，一个特定的结果可以从它现在的状态转移到任何状态，或者保持原状；
 - 每个阶段从一个状态转移到其他状态的概率用一个转移矩阵表示，矩阵每行的各元素在0到1之间，每行的和为1。
- 选举投票趋势预测
 - 以美国大选为例，首先取得过去十次选举的历史数据，然后根据历史数据得到选民意向的转移矩阵，转移矩阵如下

下一状态 当前状态	共和党	民主党	独立候选人
共和党	0.75	0.05	0.20
民主党	0.20	0.60	0.20
独立候选人	0.40	0.20	0.40

- 比如，当前状态的共和党转移到下一状态的共和党的概率是0.75，以此类推如下关系：



- 然后我们可以构造出差分表达式(共和党R，民主党D，独立候选人I)：
也就是下个状态等于前一个状态的所有可能之和

$$\begin{aligned} R_{n+1} &= 0.75R_n + 0.20D_n + 0.40I_n \\ D_{n+1} &= 0.05R_n + 0.60D_n + 0.20I_n \\ I_{n+1} &= 0.20R_n + 0.20D_n + 0.40I_n \end{aligned}$$

- 通过求解差分方程组，预测出选民投票意向的长期趋势

- plt.annotate是标记文本，如

plt.annotate('DemocraticParty', xy=(5, 0.2)) 中，xy=(a, b)是文字的位置，需要自己多试几次调一下

```

1  import matplotlib.pyplot as plt
2
3  if __name__ == '__main__':
4
5      RLIST = [1 / 3]
6      DLIST = [1 / 3]
7      ILIST = [1 / 3]
8      for i in range(40):
9          R = RLIST[i] * 0.75 + DLIST[i] * 0.20 + ILIST[i] * 0.40
10         RLIST.append(R)
11         D = RLIST[i] * 0.05 + DLIST[i] * 0.60 + ILIST[i] * 0.20
12         DLIST.append(D)
13         I = RLIST[i] * 0.20 + DLIST[i] * 0.20 + ILIST[i] * 0.40
14         ILIST.append(I)
15         plt.plot(RLIST)
16         plt.plot(DLIST)
17         plt.plot(ILIST)
18         plt.xlabel('Time')

```

```

19     plt.ylabel('Voting percent')
20     plt.annotate('DemocraticParty', xy=(5, 0.2))
21     plt.annotate('RepublicanParty', xy=(5, 0.5))
22     plt.annotate('IndependentCandidate', xy=(5, 0.25))
23     plt.show()
24     print(RLIST, DLIST, ILIST)
25
26     print('预测的最后一年: RLIST: {}, DLIST: {}, ILIST: {}'.format(RLIST[-1],
27                               DLIST[-1], ILIST[-1]))

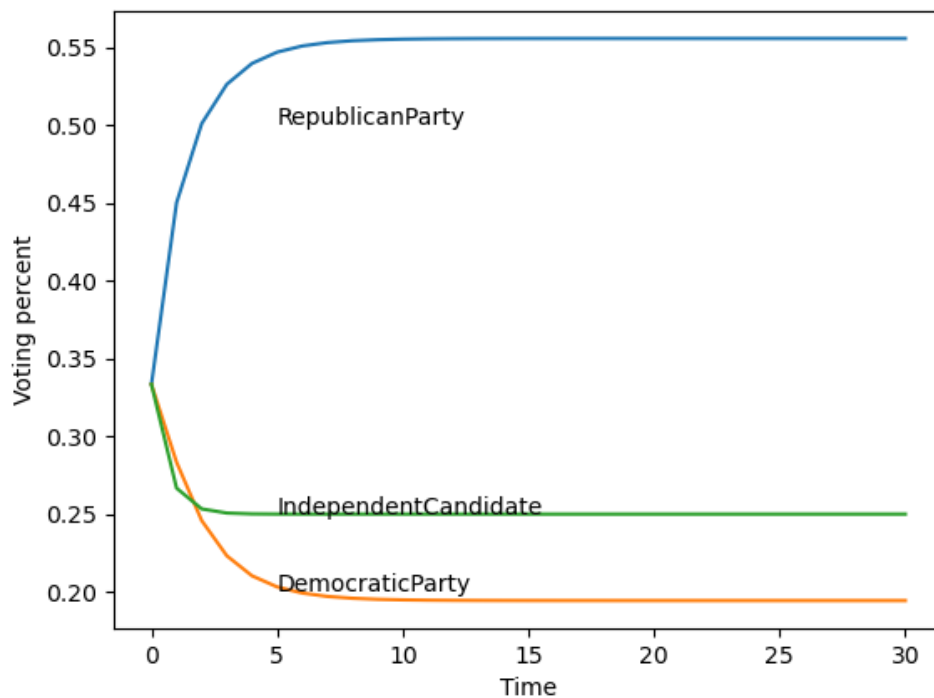
```

- 遍历画出每一年，这是最后一年的图

```

1  .....
2  .....
3  预测的最后一年: RLIST: 0.55555555555483689, DLIST: 0.1944444444516318, ILIST:
   0.25000000000000002
4
5  Process finished with exit code 0

```



- 最后得到的长期趋势是：
 - 56%的人选共和党、
 - 19%的人选民主党、
 - 25%的人选独立候选人

灰色与模糊

多层模糊评价

- 见相关资料

模糊c均值聚类

- 原理见资料
- 聚类的算法实现参考我之前sklearn里面的聚类算法实现就行了

灰色预测(经典常用)

- 灰色预测是用灰色模式GM(1,1)来进行定量分析的，通常分为以下几类：
 - 灰色时间序列预测。用等时距观测到的反映预测对象特征的一系列数量（如产量、销量、人口数量、存款数量、利率等）构造灰色预测模型，预测未来某一时刻的特征量，或者达到某特征量的时间。
 - 畸变预测（灾变预测）。通过模型预测异常值出现的时刻，预测异常值什么时候出现在特定时区内。
 - 波形预测，或拓扑预测，通过灰色模型预测事物未来变动的轨迹。
 - 系统预测，对系统行为特征指标建立一族相互关联的灰色预测理论模型，在预测系统整体变化的同时，预测系统各个环节的变化。
- 上述灰色预测方法的共同特征是：
 - 允许少数据预测；
 - 允许对灰因果律实践进行预测，例如：
 - 灰因白果律事件：粮食生产预测(就是结果产量是已知的，中间受什么因素影响是未知的)
 - 白因灰果律事件：开放项目前景预测(过程已知，但是结果前景未知)
 - 具有可检验性（事前检验：建模可行性级比检验；模型检验：建模精度检验；
预测检验：预测滚动检验）

- 模型理论部分见资料或者网上看

算法步骤

- 要使用灰色预测模型，首先看看适不适用
- 如果级比都落在可容覆盖范围内，就直接用
否则做平移变换

(1) 数据的级比检验

为了保证灰色预测的可行性，需要对原始序列数据进行级比检验。
对原始数据列

$$X^{(0)} = (x^{(0)}(1), x^{(0)}(2), \dots, x^{(0)}(n))$$

计算序列的级比：

$$\lambda(k) = \frac{x^{(0)}(k-1)}{x^{(0)}(k)}, k = 2, \dots, n$$

若所有的级比 $\lambda(k)$ 都落在可容覆盖 $\Theta = (e^{-2/(n+1)}, e^{2/(n+2)})$ 内，则可进行灰色预测；
否则需要对 $X^{(0)}$ 做平移变换， $Y^{(0)} = X^{(0)} + c$ ，使得 $Y^{(0)}$ 满足级比要求。

(2) 建立GM(1,1)模型，计算出预测值列。

(3) 检验预测值：

①相对残差检验，计算

$$\varepsilon(k) = \frac{x^{(0)}(k+1)}{x^{(0)}(k)}, k = 2, \dots, n$$

若 $\varepsilon(k) < 0.2$ ，则认为达到一般要求，若 $\varepsilon(k) < 0.1$ ，则认为达到较高要求；

②级比偏差值检验

根据前面计算出来的级比 和发展系数，计算相应的级比偏差：

若，则认为达到一般要求，若，则认为达到较高要求。

(4) 利用模型进行预测

- 代码部分是可以使用cuda加速的
- 只需改87行的输入数据和93行的预测个数m

```
1 import torch as th
2 import numpy as np
3
4
5 class GM:
6     def __init__(self):
7         # 判断是否可用 gpu 编程，大量级计算使用GPU
8         self._is_gpu = False # th.cuda.is_available()
9
10    def fit(self, dt: list or np.ndarray):
11        self._df: th.Tensor = th.from_numpy(np.array(dt, dtype=np.float32))
12        if self._is_gpu:
13            self._df.cuda()
14        self._n: int = len(self._df)
15        self._x, self._max_value = self._sigmoid(self._df)
16        z: th.Tensor = self._next_to_mean(th.cumsum(self._x, dim=0))
17        self.coef: th.Tensor = self._coefficient(self._x, z)
18        del z
19        self._x0: th.Tensor = self._x[0]
20        self._pre: th.Tensor = self._pred()
21
```

```

22 # 归一化
23 def _sigmod(self, x: th.Tensor):
24     _maxv: th.Tensor = th.max(x)
25     return th.div(x, _maxv), _maxv
26
27 # 计算紧邻均值数列
28 def _next_to_mean(self, x_1: th.Tensor):
29     z: th.Tensor = th.zeros(self._n - 1)
30     if self._is_gpu:
31         z.cuda()
32     for i in range(1, self._n): # 下标从0开始, 取不到最大值
33         z[i - 1] = 0.5 * x_1[i] + 0.5 * x_1[i - 1]
34     return z
35
36 # 计算系数 a,b
37 def _coefficient(self, x: th.Tensor, z: th.Tensor):
38     B: th.Tensor = th.stack((-1 * z, th.ones(self._n - 1)), dim=1)
39     Y: th.Tensor = th.tensor(x[1:], dtype=th.float32).reshape((-1, 1))
40     if self._is_gpu:
41         B.cuda()
42         Y.cuda()
43
44     # 返回的是a和b的向量转置, 第一个是a 第二个是b;
45     return th.matmul(th.matmul(th.inverse(th.matmul(B.t(), B)), B.t()), Y)
46
47 def _pred(self, start: int = 1, end: int = 0):
48     les: int = self._n + end
49     resut: th.Tensor = th.zeros(les)
50
51     if self._is_gpu:
52         resut.cuda()
53     resut[0] = self._x0
54     for i in range(start, les):
55         resut[i] = (self._x0 - (self.coef[1] / self.coef[0])) * \
56             (1 - th.exp(self.coef[0])) * th.exp(-1 * self.coef[0] * (i))
57     del les
58     return resut
59
60 # 计算绝对误差
61 def confidence(self):
62     return round((th.sum(th.abs(th.div((self._x - self._pre), self._x))) /
self._n).item(), 4)
63
64 # 预测个数, 默认个数大于等于0,
65 def predict(self, m: int = 1, decimals: int = 4):
66     y_pred: th.Tensor = th.mul(self._pre, self._max_value)
67     y_pred_ = th.zeros(1)
68     if m < 0:
69         return "预测个数需大于等于0"
70     elif m > 0:

```

```

71         y_pred_: th.Tensor = self._pred(self._n, m)[-m:].mul(self._max_value)
72     else:
73         if self._is_gpu:
74             return list(map(lambda _: round(_, decimals),
y_pred.cpu().numpy().tolist()))
75         else:
76             return list(map(lambda _: round(_, decimals),
y_pred.numpy().tolist()))
77
78     # cat 拼接 0 x水平拼接, 1y垂直拼接
79     result: th.Tensor = th.cat((y_pred, y_pred_), dim=0)
80     del y_pred, y_pred_
81     if self._is_gpu:
82         return list(map(lambda _: round(_, decimals),
result.cpu().numpy().tolist()))
83     return list(map(lambda _: round(_, decimals), result.numpy().tolist()))
84
85
86 if __name__ == "__main__":
87     ls = np.arange(91, 100, 2) # ls是原始的值
88     print(type(ls))
89     gm = GM()
90     gm.fit(ls)
91     print('绝对误差: ', gm.confidence())
92     print('原始值: ', ls)
93     print('预测: ', gm.predict(m=2)) # m是2代表要预测后面两个值
94

```

```

1 <class 'numpy.ndarray'>
2 绝对误差: 0.0002
3 原始值: [91 93 95 97 99]
4 预测: [91.0, 93.0178, 94.9758, 96.9751, 99.0164, 101.1007, 103.2289]
5
6 Process finished with exit code 0

```

蒙特卡罗

蒙特卡罗算法

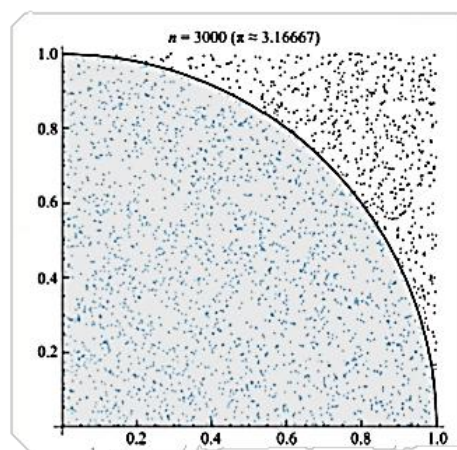
- 由冯·诺依曼提出来的
- 蒙特·卡罗(Monte Carlo method)又称统计模拟方法，一种以概率统计理论为指导的数值计算方法。是指使用随机数(或者伪随机数)来解决很多计算问题的方法。
- 基本思想

当所求解问题是某种随机事件出现的概率，或者是某个随机变量的期望值时，通过某种“实验”的方法，以这种事件出现的频率估计这一随机事件的概率，或者得到这个随机变量的某些数字特征，并将其作为问题的解。

举例

- 蒙特卡罗方法求圆周率

基本思想：在图中区域产生足够多的随机数点，然后计算落在圆内的点的个数与总个数的比值再乘以4，就是圆周率。



- `random.random()`得到的是0到1之间的(伪)随机数，若要某个整数范围，则用`random.randint(a, b)`，是有包括a和b的

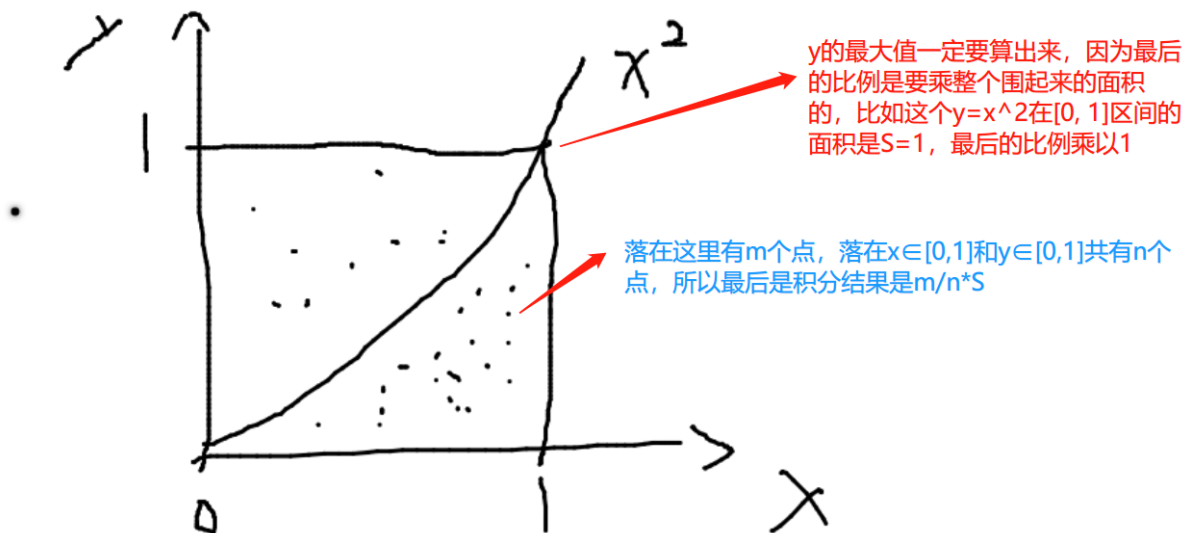
```
1 import math
2 import random
3
4 if __name__ == '__main__':
5
6     M = input('请输入一个较大的整数')
7     N = 0
8     for i in range(int(M)):
9         x = random.random()
10        y = random.random()
11        if math.sqrt(x ** 2 + y ** 2) < 1:
12            N += 1
13            pi = 4 * N / int(M)
14            # print(pi)
15    print(pi)
16
```

```
1 请输入一个较大的整数9999999
2 3.1417339141733915
3
4 Process finished with exit code 0
```

- 蒙特卡罗求定积分

利用python计算函数 $y=x^2$ 在 $[0,1]$ 区间的定积分

基本思想：和上例相似



```
1  """
2  求f(x) = x^2的定积分
3  """
4
5  import random
6
7  if __name__ == '__main__':
8      n = int(input('请输入一个较大的整数:'))
9      m = 0
10     for i in range(n):
11         x = random.random()
12         y = random.random()
13         if y < x ** 2: # 找到落在f(x)下面的点
14             m += 1
15     R = m / n # 这里因为总面积是1所以省略乘以1了
16     print(R)
17
```

```
1  请输入一个较大的整数:999999
2  0.33292533292533294
3
4  Process finished with exit code 0
```

三门问题

- 背景

三门问题 (Monty Hall problem) 亦称为蒙提霍尔问题，出自美国电视游戏节目Let's Make a Deal。参赛者会看见三扇关闭了的门，其中一扇的后面有一辆汽车，选中后面有车的那扇门可赢得该汽车，另外两扇门则各藏有一只山羊。当参赛者选定了一扇门，但未去开启它的时候，节目主持人开启剩下两扇门的其中一扇，露出其中一只山羊。主持人其后问参赛者要不要换另一扇仍然关上的门。问题是：换另一扇门是否会增加参赛者赢得汽车的几率？如果严格按照上述条件，即主持人清楚地知道，自己打开的那扇门后面是羊，那么答案是会。不换门的话，赢得汽车的几率是1/3，换门的话，赢得汽车的几率是2/3

- 应用蒙特卡罗重点在使用随机数来模拟类似于赌博问题的赢率问题，通过多次模拟得到所要计算值的模拟值
- 解决思路：

在三门问题中，用0、1、2分代表三扇门的编号，在 $[0, 2]$ 之间随机生成一个整数代表奖品所在门的编号prize，再次在 $[0, 2]$ 之间随机生成一个整数代表参赛者所选择的门的编号guess。用变量change代表游戏中的换门(true)与不换门(false)

主持人一定会开启一扇没有奖品的门，假设一开始就猜中，剩下最后一扇肯定不中奖，否则，最后一扇肯定中奖



https://blog.csdn.net/beidu_41327283

```
1 import random
2
3
4 def play(change):
5     prize = random.randint(0, 2)
6     guess = random.randint(0, 2)
7     if prize == guess:
8         if change:
9             return False
10        else:
11            return True
12    else:
13        if change:
14            return True
15        else:
16            return False
17
18
19 def winRate(change, N):
20     win = 0
21     for i in range(0, N):
22         if play(change):
23             win = win + 1
```

```

24         # print('中奖率为: ')
25         # print(win / N)
26     print('中奖率为: ')
27     print(win / N)
28
29
30 if __name__ == '__main__':
31     N = 100000
32     print('玩' + str(N) + '次, 每一次都换门:')
33     winRate(True, N)
34     print()
35     print('玩' + str(N) + '次, 每一次都不换门:')
36     winRate(False, N)
37

```

```

1 玩100000次, 每一次都换门:
2 中奖率为:
3 0.66748
4
5 玩100000次, 每一次都不换门:
6 中奖率为:
7 0.33008
8
9 Process finished with exit code 0

```

巧克力豆问题

- 见相关资料

时间序列

- 时序问题也可使用神经网络里面的LSTM(长短时记忆)
- 预测的是近期的, 不是预测长远的(长远预测需要挖掘更多特征, 深度学习那块的)
- 均方差:

```

1 from sklearn.metrics import mean_squared_error
2 from math import sqrt
3 rms = sqrt(mean_squared_error(test, pred))# 把实际和预测的放进去
4 print(rms)

```

- 见<https://www.bilibili.com/video/BV12h411d7Dm?p=10>

简单指数平滑法

- 导入数据
- 切分数据
- 代码适当修改和测试

```

1  from statsmodels.tsa.api import SimpleExpSmoothing
2  import numpy as np
3  import pandas as pd
4  import matplotlib.pyplot as plt
5  from sklearn.metrics import mean_squared_error
6  from math import sqrt
7
8  if __name__ == '__main__':
9      df = pd.read_csv('train.csv')
10     train = df[0:10392]
11     test = df[10392:]
12     pred = test.copy()
13     fit = SimpleExpSmoothing(np.asarray(train['列名1'])).fit(smoothing_level=0.6,
optimized=False)
14     pred['列名2'] = fit.forecast(len(test)) # 需要预测多长
15
16     # 画出来
17     plt.figure(figsize=(16, 8))
18     plt.plot(train['列名1'], label='Train')
19     plt.plot(test['列名1'], label='Test')
20     plt.plot(pred['列名2'], label='列名2')
21     plt.legend(loc='best')
22     plt.show()
23
24     # 评估
25     rms = sqrt(mean_squared_error(test['列名1'], pred))
26     print(rms)
27

```

霍尔特线性趋势法

- 考虑到数据集变化趋势的方法就叫做霍尔特线性趋势法，任何呈现某种趋势(比如商上升趋势)的数据集都可以用霍尔特线性趋势法用于预测
- 先看看是否呈现某种趋势

```

1  import statsmodels.api as sm
2  result = sm.tsa.stattools.adfuller(train['xxx'])
3  plt.show()

```

- 如果是的话就可以使用
- 代码相对上面只需改动这两句，其他适当改即可

```

1  from statsmodels.tsa.api import Holt
2  fit = Holt(np.asarray(train['xxx'])).fit(smoothing_level=0.3, smoothing_slope=0.1)

```

Holt-Winters季节性预测模型

- 体现在季节性
 - 比如一个水果店的销售情况，在夏季收入远高于其他季节
- 只需改变一点代码，选择了seasonal_period=7作为每周重复数据

```
1 from statsmodels.tsa.api import ExponentialSmoothing
2 fit1 = ExponentialSmoothing(np.asarray(train['xxx']), seasonal_periods=7,
   trend='add', seasonal='add', ).fit()
```

自回归移动平均模型 (ARIMA)

- 指数平滑模型都是基于数据中的趋势和季节性的描述，而自回归移动平均模型的目标是描述数据中彼此之间的关系。ARIMA的一个优化版就是季节性ARIMA。它像Holt-Winters季节性预测模型一样，也把数据集的季节性考虑在内。

```
1 import statsmodels.api as sm
2 pred = test.copy()
3 fit1 = sm.tsa.statespace.SARIMAX(train.xxx, order=(2, 1, 4), seasonal_order=(0, 1, 1,
   7)).fit() pred['SARIMA']=fit1.predict(start="20xx-xx-xx",end="20xx-xx-
   xx",dynamic=True)
```

SVM

```
1 clf = svm.SVC(C=0.8, kernel='rbf', gamma=20, decision_function_shape='ovr')
2 clf.fit(x_train, y_train.ravel())
```

- 调参，参数见资料
- 参考：<https://www.bilibili.com/video/BV12h411d7Dm?p=12>

svm处理非线性问题，低维映射到高维

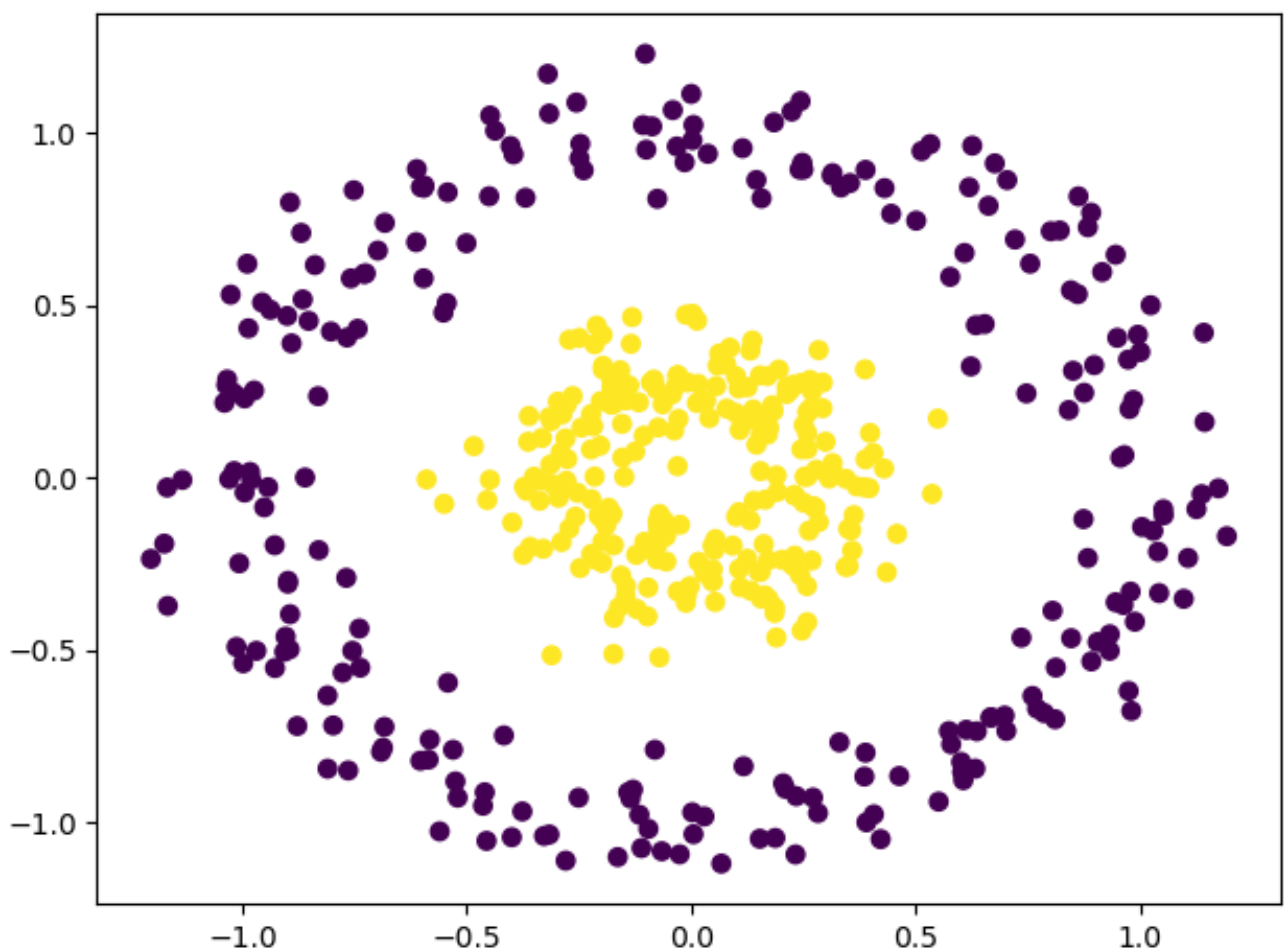
- 如2维转3维，找到切平面，在投影到2维平面，可能是个圆或者椭圆的分界线

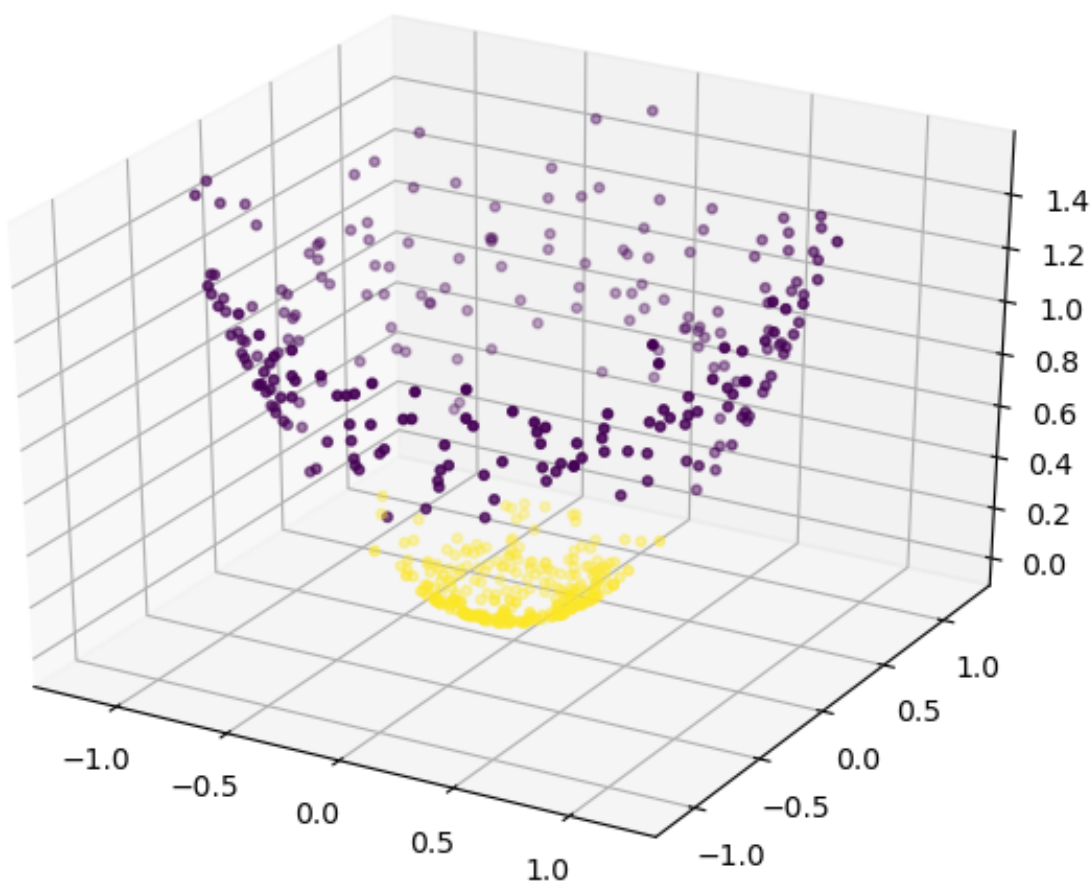
```
1 """
2 # @Time      : 2020/8/22
3 # @Author    : Jimou Chen
4 """
5 import matplotlib.pyplot as plt
6 from sklearn import datasets
```

```

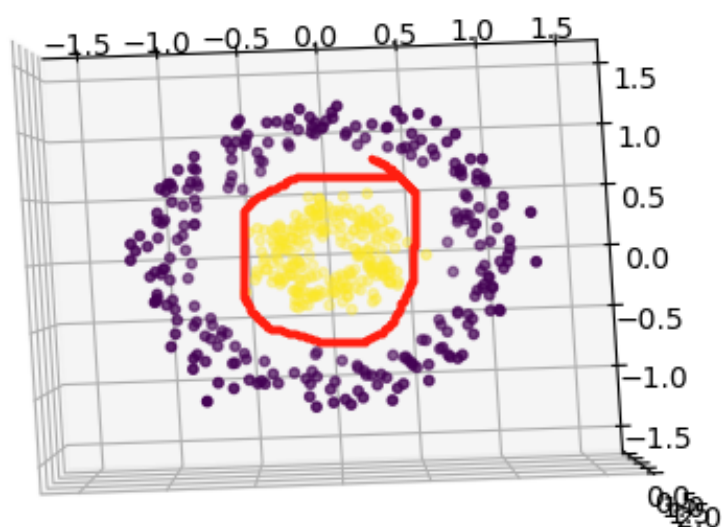
7
8 '''用于解决非线性问题'''
9
10 # 制造数据
11 x_data, y_data = datasets.make_circles(n_samples=500, factor=0.3, noise=0.1)
12 # 画出来看看
13 plt.scatter(x_data[:, 0], x_data[:, 1], c=y_data)
14 plt.show()
15
16 '''接下来把2维映射到3维'''
17
18 z_data = x_data[:, 0] ** 2 + x_data[:, 1] ** 2
19 # 画3d图
20 ax = plt.figure().add_subplot(111, projection='3d')
21 ax.scatter(x_data[:, 0], x_data[:, 1], z_data, c=y_data, s=10) # s是大小
22 plt.show()
23

```





投影到二维平面



核函数

我们可以构造核函数使得运算结果等同于非线性映射，同时运算量要远远小于非线性映射。

$$K(X_i, X_j) = \phi(X_i) \cdot \phi(X_j)$$

$$h \text{ 次多项式核函数: } K(X_i, X_j) = (X_i \cdot X_j + 1)^h$$

$$\text{高斯径向基函数核函数: } K(X_i, X_j) = e^{-\|X_i - X_j\|^2 / 2\sigma^2}$$

$$\text{S 型核函数: } K(X_i, X_j) = \tanh(\kappa X_i \cdot X_j - \delta)$$

svm推导过程

- 推导过程理解起来较为复杂，见其他机器学习教材或者相关资料

svm处理多分类问题

- 一般svm是处理二值分类问题的，如果处理多个类别的，使用下面三种方法

```
1 model = svm.SVC(decision_function_shape='ovo')
2 model = svm.SVC(decision_function_shape='ovr')
3 model = svm.SVC(probability=True)
```

- 以下是先用pca降维后，再用svm进行分类的例子

```
1 """
2 # @Time      : 2020/8/30
3 # @Author    : Jimou Chen
4 """
5 from sklearn.decomposition import PCA
6 from sklearn.model_selection import train_test_split
7 import matplotlib.pyplot as plt
8 import pandas as pd
9 from sklearn import svm
10 from sklearn.metrics import classification_report
11
12
13 data = pd.read_csv('data/wine.csv')
14 y_data = data.iloc[:, 0]
15 x_data = data.iloc[:, 1:]
16
17 x_train, x_test, y_train, y_test = train_test_split(x_data, y_data)
18 pca = PCA(n_components=2)
```

```

19 new_data = pca.fit_transform(x_data)
20
21 # 画出来看一下
22 plt.scatter(new_data[:, 0], new_data[:, 1], c=y_data)
23 plt.show()
24
25 # 建模预测, 多分类的三种方法, 有时候会警告, 不影响
26 # model = svm.SVC(decision_function_shape='ovo')
27 model = svm.SVC(decision_function_shape='ovr')
28 # model = svm.SVC(probability=True)
29 model.fit(x_train, y_train)
30
31 prediction = model.predict(x_data)
32
33 print(model.score(x_test, y_test))
34 print(classification_report(y_data, prediction))
35
36 # 画出预测的
37 plt.scatter(new_data[:, 0], new_data[:, 1], c=prediction)
38 plt.show()

```

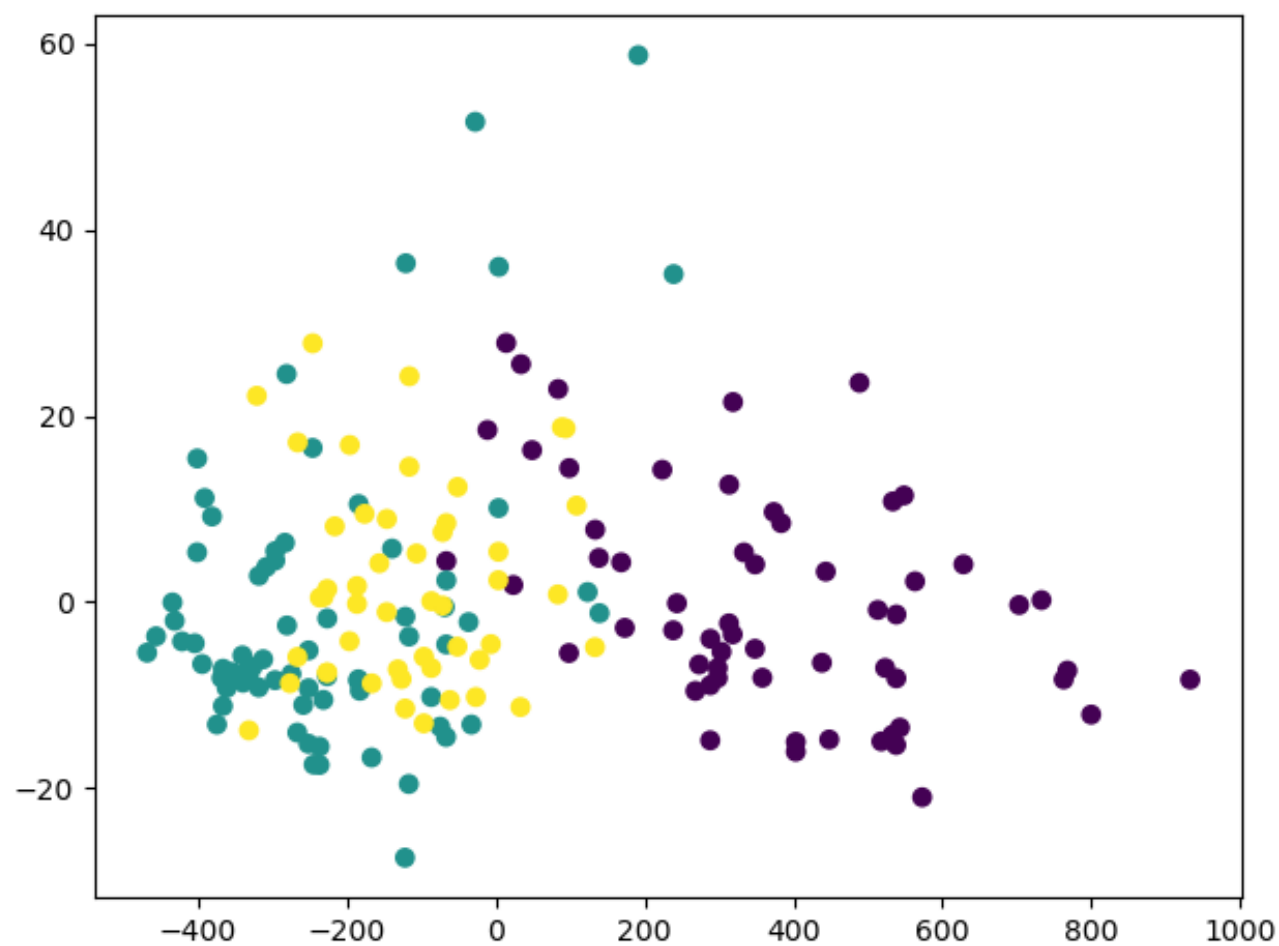
- 结果:

```

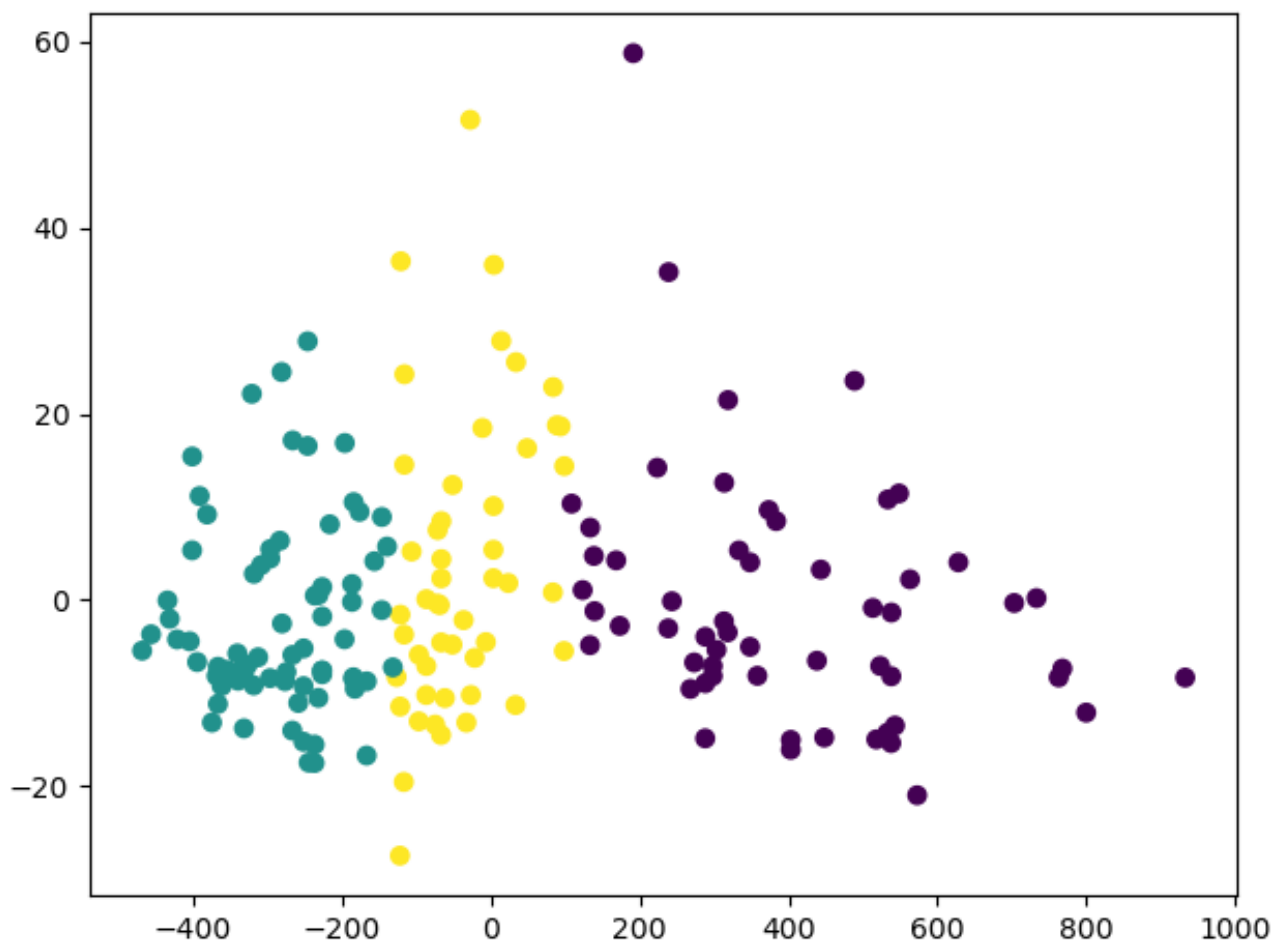
1 0.6888888888888889
2           precision    recall  f1-score   support
3
4      1         0.89      0.85      0.87         59
5      2         0.70      0.72      0.71         71
6      3         0.49      0.50      0.49         48
7
8   accuracy                        0.70         178
9   macro avg         0.69      0.69      0.69         178
10  weighted avg         0.71      0.70      0.70         178
11
12
13 Process finished with exit code 0

```

- 原始的



- 预测的



- 改进

上面这个题用随机森林的效果会更好，无论降到2/3/5维，预测效果几乎完美

```
model = RandomForestClassifier(n_estimators=100)
```

聚类

K-Means模型

- 算法思想：以空间中k个点为中心进行聚类，对最靠近他们的对象归类。通过迭代的方法，逐次更新各聚类中心的值，直至得到最好的聚类结果
- 跟分类相比，没有给定已知标签
- 用sklearn自带的KMeans效果要比自己实现的Kmeans效果好
- 用sklearn实现如下，以下是以二维为例子
 - 修改第6行的文件文件即可，kmeans.txt的格式见：
 - <https://www.lanzoui.com/iykw2re7qxc>
密码:hmse
 - 这里的文件可以灵活处理，如果是csv格式的需要稍微处理

```
1 import numpy as np
```

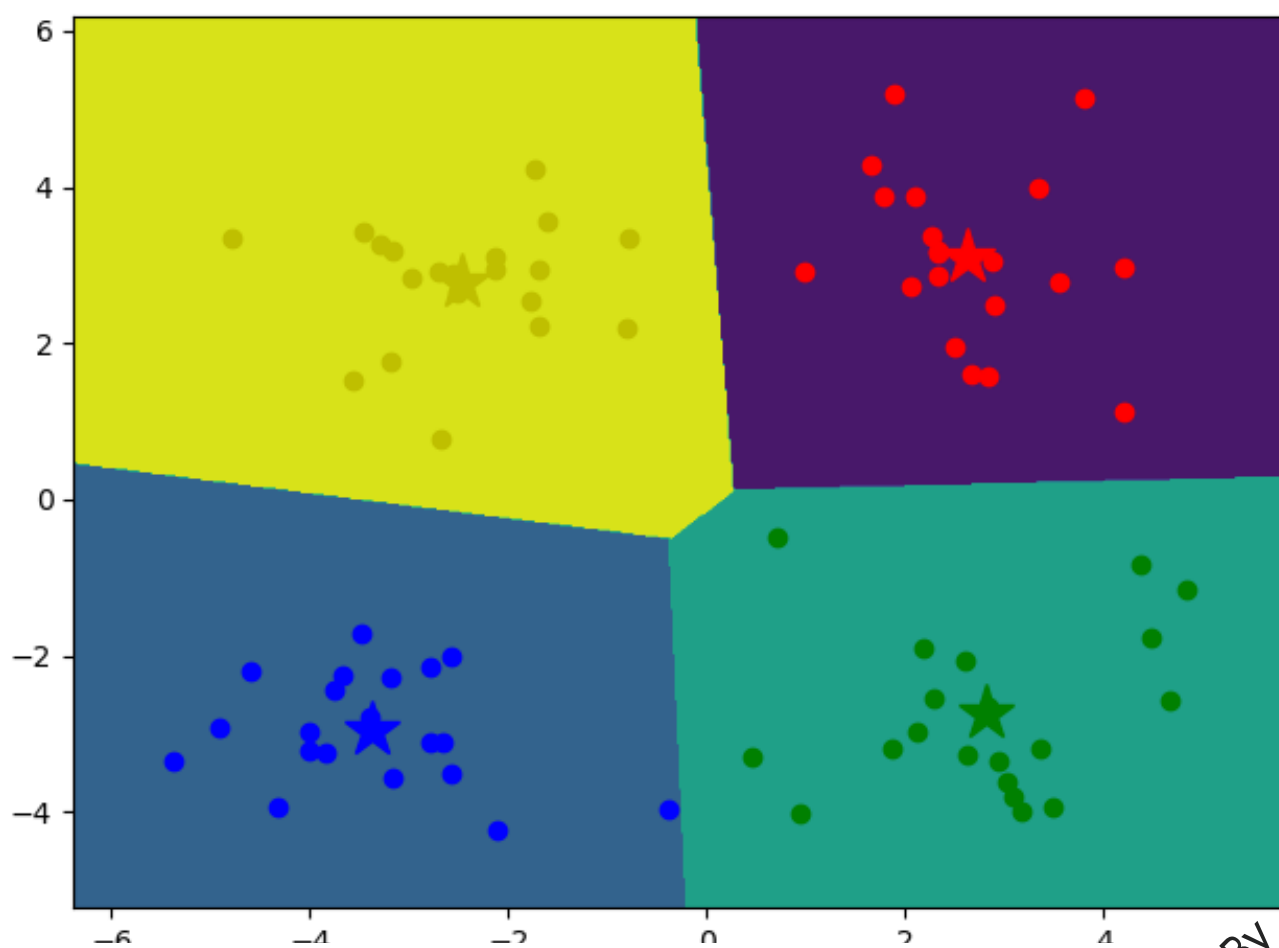
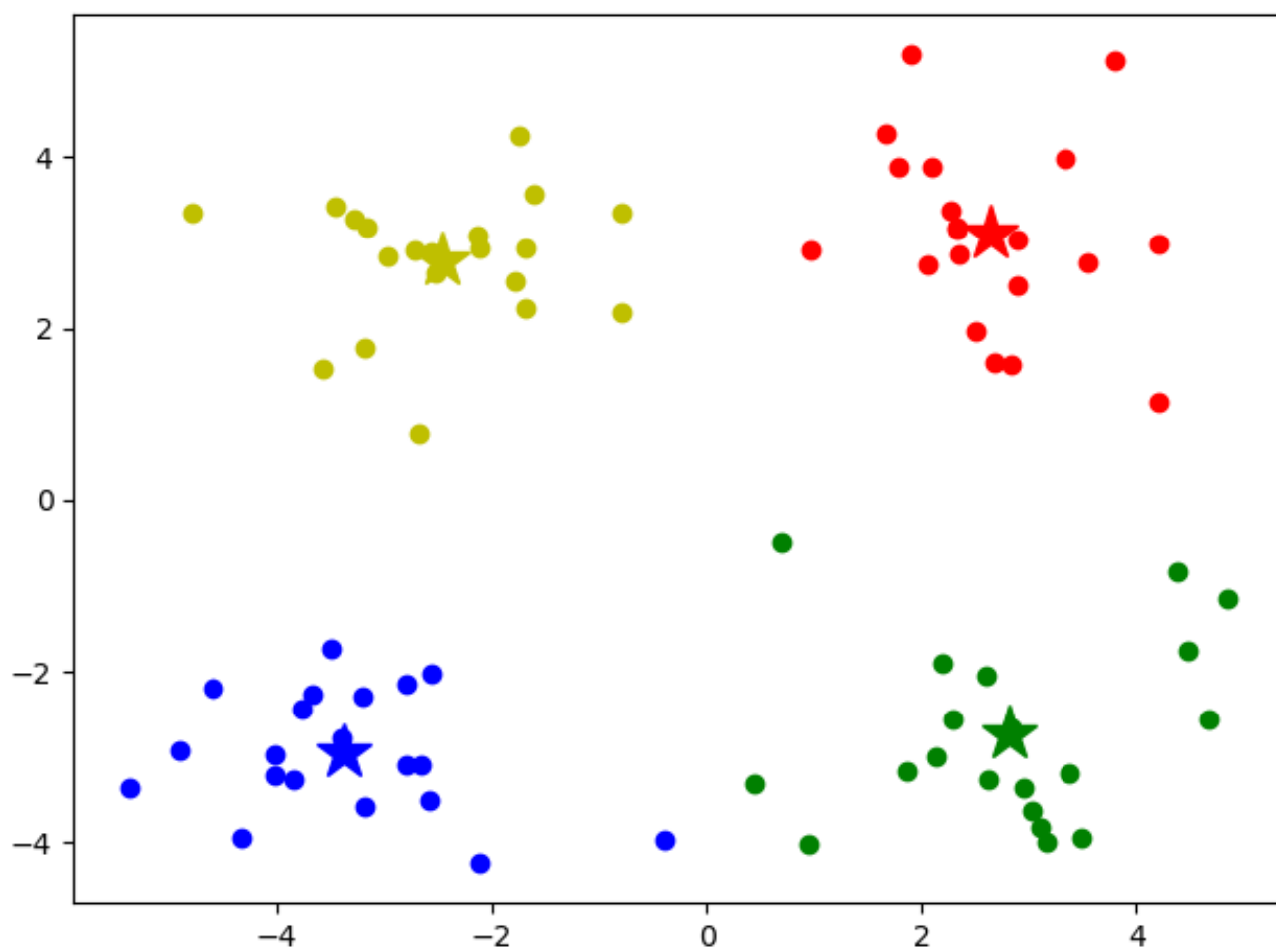
```

2 import matplotlib.pyplot as plt
3 from sklearn.cluster import KMeans # 导入KMeans
4
5 # 载入数据
6 data = np.genfromtxt('kmeans.txt', delimiter=' ')
7 # 设置k值
8 k = 4
9
10 # 建模
11 model = KMeans(n_clusters=k)
12 model.fit(data)
13
14 # 各分类的中心点坐标
15 centers = model.cluster_centers_
16 print(centers)
17
18 # 预测结果
19 prediction = model.predict(data)
20 print(prediction)
21
22 # 画图
23 colors = ['or', 'ob', 'og', 'oy']
24 for i, d in enumerate(data):
25     plt.plot(d[0], d[1], colors[prediction[i]])
26
27 # 画出各个分类的中心点
28 mark = ['*r', '*b', '*g', '*y']
29 for i, center in enumerate(centers):
30     plt.plot(center[0], center[1], mark[i], markersize=20)
31
32 plt.show()
33
34 # 获取数据值所在的范围
35 x_min, x_max = data[:, 0].min() - 1, data[:, 0].max() + 1
36 y_min, y_max = data[:, 1].min() - 1, data[:, 1].max() + 1
37
38 # 生成网格矩阵
39 xx, yy = np.meshgrid(np.arange(x_min, x_max, 0.02),
40                      np.arange(y_min, y_max, 0.02))
41
42 z = model.predict(np.c_[xx.ravel(), yy.ravel()]) # ravel与flatten类似，多维数据转一维。
43 # flatten不会改变原始数据，ravel会改变原始数据
44 z = z.reshape(xx.shape)
45 # 等高线图
46 cs = plt.contourf(xx, yy, z)
47 # 显示结果
48 # 画出各个数据点，用不同颜色表示分类
49 mark = ['or', 'ob', 'og', 'oy']
50 for i, d in enumerate(data):
    plt.plot(d[0], d[1], mark[prediction[i]])

```

```
51
52 # 画出各个分类的中心点
53 mark = ['*r', '*b', '*g', '*y']
54 for i, center in enumerate(centers):
55     plt.plot(center[0], center[1], mark[i], markersize=20)
56
57 plt.show()
58
```

```
1 [[ 2.6265299  3.10868015]
2  [-3.38237045 -2.9473363 ]
3  [ 2.80293085 -2.7315146 ]
4  [-2.46154315  2.78737555]]
5 [0 3 2 1 0 3 2 1 0 3 2 1 0 3 2 1 0 3 2 1 0 3 2 1 0 3 2 1 0
6  3 2 1 0 3 2 1 0 3 2 1 0 3 2 1 0 3 2 1 0 3 2 1 0 3 2 1 0 3
7  2 1 0 3 2 1]
8
9 Process finished with exit code 0
```



model.fit(data) 只传一个

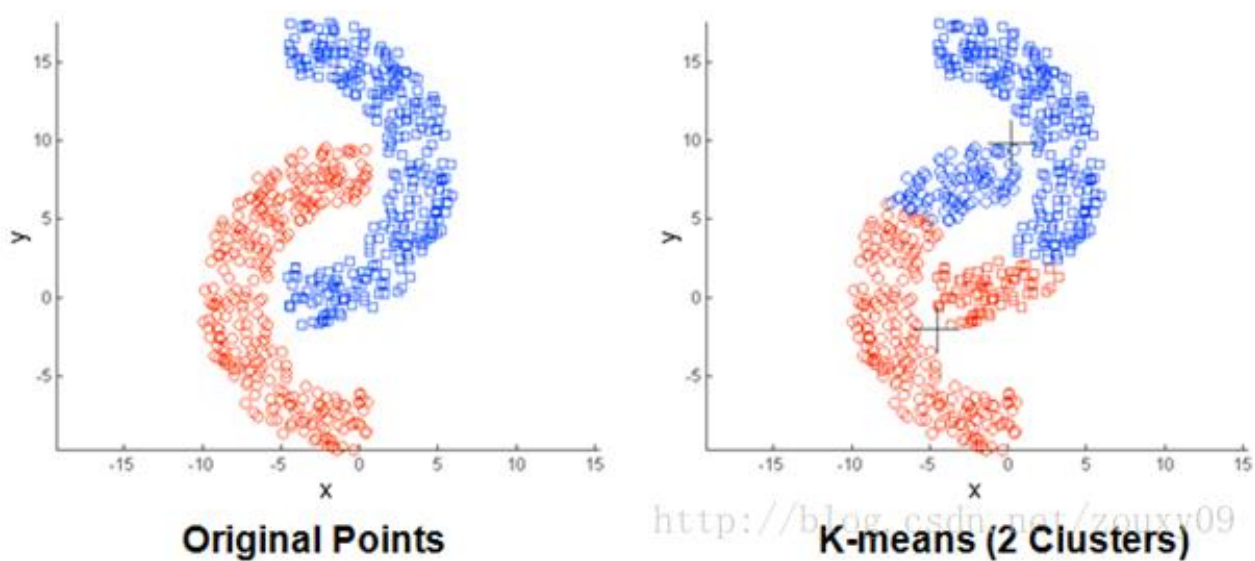
MiniBatchKMeans模型

- Mini Batch K-Means算法是K-Means算法的变种
- 与K均值算法相比，数据的更新是在每一个小的样本集上。Mini Batch K-Means比K-Means有更快的收敛速度，但同时也降低了聚类效果，但是在实际项目中却表现得不明显

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.cluster import MiniBatchKMeans
4
5 '''适用于数据多的情况，但是一般情况还是用KMeans'''
6
7 data = np.genfromtxt('kmeans.txt', delimiter=' ')
8 k = 4
9
10 model = MiniBatchKMeans(n_clusters=k)
11 model.fit(data)
12
13 # 分类中心坐标
14 centers = model.cluster_centers_
15 print(centers)
16
17 # 预测结果
18 pred_res = model.predict(data)
19 print(pred_res)
20
21 # 画图
22 colors = ['or', 'ob', 'og', 'oy']
23 for i, d in enumerate(data):
24     plt.plot(d[0], d[1], colors[pred_res[i]])
25
26 # 画出各个分类的中心点
27 mark = ['*r', '*b', '*g', '*y']
28 for i, center in enumerate(centers):
29     plt.plot(center[0], center[1], mark[i], markersize=20)
30
31 plt.show()
32
```

- KMeans的4个问题，可参考解释：<https://www.bilibili.com/video/BV1Rt411q7WJ?p=65>
 - 对k个初始质心的选择比较敏感，容易陷入局部最小值
 - k值的选择是用户指定的，不同的k得到的结果会有挺大的不同

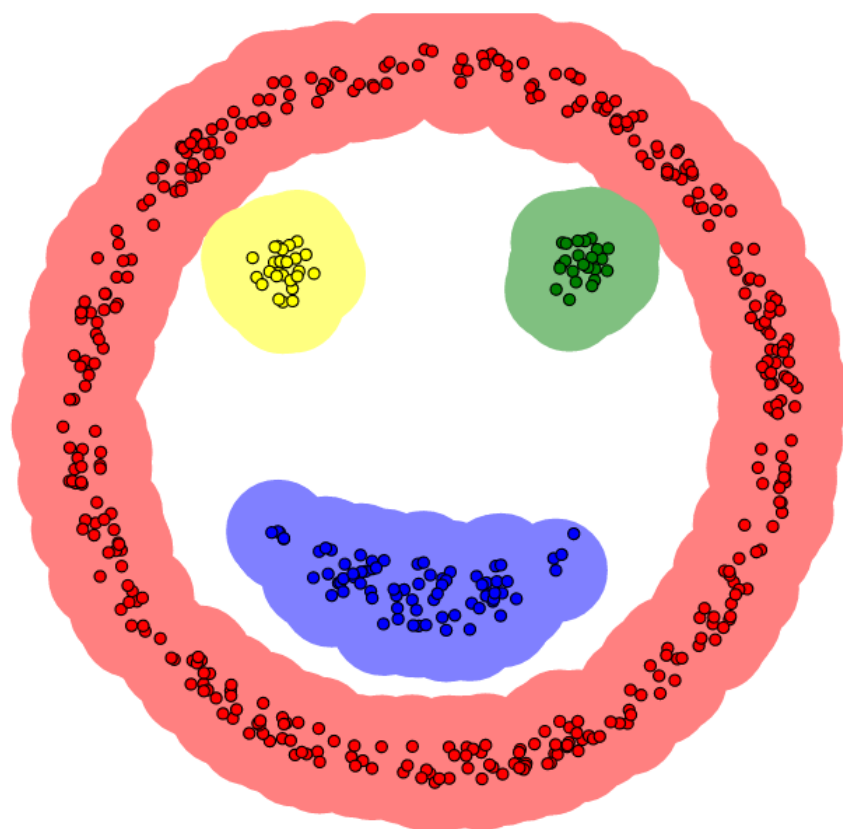
- 存在局限性，如下面这种非球状的数据分布就搞不定了



- 数据比较大的时候，收敛会比较慢

DBSCAN模型

- 本算法将具有足够高密度的区域划分为簇，并可以发现任何形状的聚类

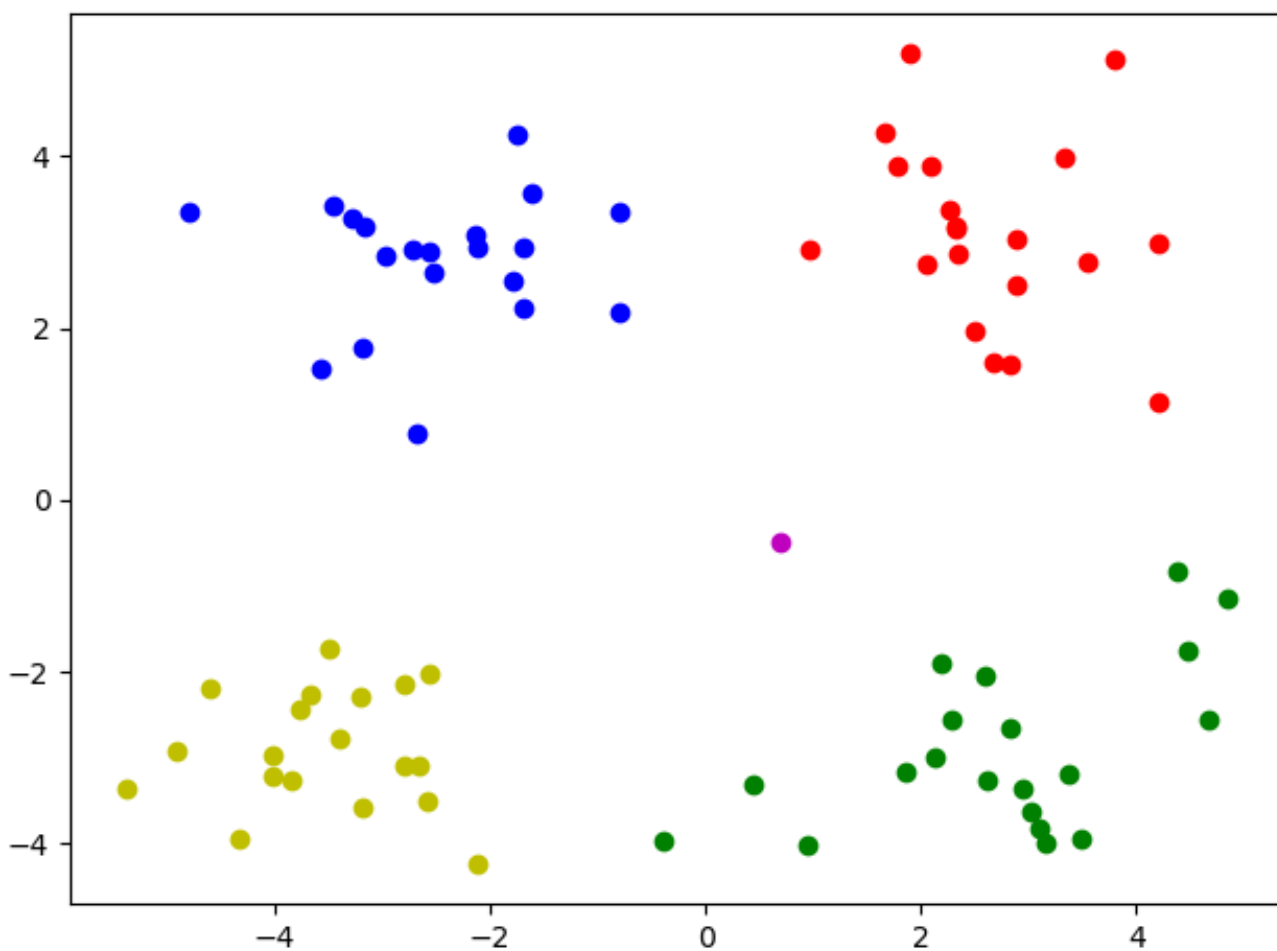


- 算法思想: <https://www.bilibili.com/video/BV1Rt411q7WJ?p=69>
- 缺点:
 - 当数据量增大时，要求较大的内存支持I/O消耗也很大。

- 当空间聚类的密度不均匀、聚类间距差相差很大时，聚类质量较差。
- DBSCAN和K-MEANS比较：
 - DBSCAN不需要输入聚类个数。
 - 聚类簇的形状没有要求。
 - 可以在需要时输入过滤噪声的参数。
- 关键参数
 - `eps`
 - `min_point`

即通过调整`eps=`, `min_samples=`来找到一个最好的效果

```
1 from sklearn.cluster import DBSCAN
2 import numpy as np
3 import matplotlib.pyplot as plt
4
5 data = np.genfromtxt('kmeans.txt', delimiter=' ')
6 # 建模
7 # esp距离阈值, min_samples在esp领域里面的样本数
8 model = DBSCAN(eps=1.5, min_samples=4)
9 model.fit(data)
10
11 # 预测,fit_predict是先拟合后预测, DBSCAN没有predict方法
12 pred = model.fit_predict(data)
13 print(pred) # 预测值为-1的是噪点
14
15 # 画出各个数据点
16 mark = ['or', 'ob', 'og', 'oy', 'ok', 'om']
17 for i, d in enumerate(data):
18     plt.plot(d[0], d[1], mark[pred[i]])
19 plt.show()
20
```



1	[0	1	2	3	0	1	2	3	0	1	2	3	0	1	-1	2	0	1	2	3	0	1	2	3
2		0	1	2	3	0	1	2	3	0	1	2	3	0	1	2	3	0	1	2	3	0	1	2	3
3		0	1	2	3	0	1	2	3	0	1	2	3	0	1	2	3	0	1	2	3	0	1	2	3
4		0	1	2	3	0	1	2	3]																

- 预测值是-1 的代表噪点