

### 最小二乘法

### 耿修瑞

中国科学院空天信息创新研究院 gengxr@sina.com.cn

2025.3

### 主要内容



- □问题背景
- □基本原理
- □变量问题
- □几何解释
- □概率解释
- □约束最小二乘
- □最小二乘法的应用

### 问题背景





1766年,德国有一位名叫约翰·提丢斯(Johann Daniel Titius)的中学教师,写下了下面的数列:

$$\frac{(3 \times 2^n + 4)}{10} \qquad n = -\infty, 0, 1, 2, 3, \dots$$

令提丢斯惊奇的是,他发现这个数列的每一项与当时已知的六大行星(即水星、金星、地球、火星、木星、土星)到太阳的距离比例(地球到太阳的距离定为1个单位)有着一定的联系。

### 问题背景



提丢斯的朋友, 天文学家波得深知这一发现的重要意义, 就于1772年公布了提丢斯的这一发现,这串数从此引起 了科学家的极大重视:并被称为提丢斯——波得定律即太 阳系行星与太阳的平均距离. 当时,人们还没有发现天王 星、海王星,以为土星就是距太阳最远的行星。1781年 英籍德国人赫歇尔 在接近19.6的位置上(即数列中的 第八项)发现了天王星,从此,人们就对这一定则深信不 疑了。根据这一定则,在数列的第五项即2.8的位置上也 应该对应一颗行星, 只是还没有被发现。于是, 许多天文 学家和天文爱好者便以极大的热情, 踏上了寻找这颗新行 星的征程。

### 问题背景



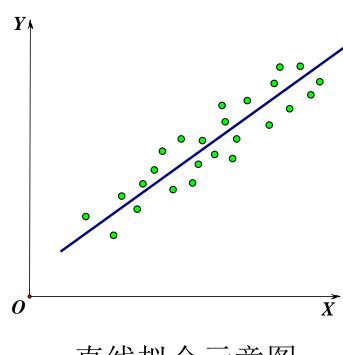
1801年, 意大利天文学家朱赛普-皮亚齐发现了第 一颗小行星谷神星。经过40天的跟踪观测后,由于 谷神星运行至太阳背后, 使得皮亚齐失去了谷神星 的位置。随后全世界的科学家利用皮亚齐的观测数 据开始寻找谷神星,但是根据大多数人计算的结果 来寻找谷神星都没有结果。时年24岁的高斯也采用 了一种新方法(即为最小二乘法)计算了谷神星的 轨道。奥地利天文学家海因里希·奥尔伯斯根据高斯 计算出来的轨道重新发现了谷神星。高斯使用的最 小二乘法的方法发表于1809年他的著作《天体运动 论》中,而法国科学家<u>勒让德</u>于1806年独立发现最 小二乘法, 但因不为时人所知而默默无闻。两人曾 为谁最早创立最小二乘法原理发生争执。最小二乘 法自创立以来,在自然科学乃至社会科学的各个领 域产生了广泛的应用。



高 斯 ( Gauss, 1777~1855), 德国数 学家,最小二乘法创始人。



最小二乘法通常用来研究两个变量或者多个变量之间的关系。以常用的直线回归为例,已知一组观测点 (x<sub>1</sub>,y<sub>1</sub>),(x<sub>2</sub>,y<sub>2</sub>),…,(x<sub>n</sub>,y<sub>n</sub>)分布在直角坐标系中(如图),如何用一条直线最佳的拟合这些散点?

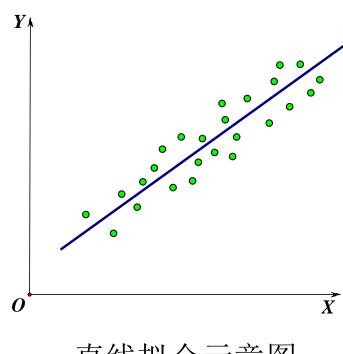


直线拟合示意图

$$y = ax + b$$



最小二乘法通常用来研究两个变量或者多个变量之间的关系。以常用的直线回归为例,已知一组观测点 (x<sub>1</sub>,y<sub>1</sub>),(x<sub>2</sub>,y<sub>2</sub>),…,(x<sub>n</sub>,y<sub>n</sub>)分布在直角坐标系中(如图),如何用一条直线最佳的拟合这些散点?



直线拟合示意图

$$y = ax + b$$



### 最小二乘法的一般解法(行空间角度)

最小二乘法求解,相当于寻找a,b。使得如下公式尽量成立

$$\begin{cases} y_1 = ax_1 + b \\ y_2 = ax_2 + b \\ \vdots \\ y_n = ax_n + b \end{cases}$$

可建立优化模型如下:

$$\min_{a,b} f(a,b) = \sum_{i=1}^{n} (ax_i + b - y_i)^2$$



对自变量a,b求导,并令其为0:

$$\begin{cases} \frac{\partial f}{\partial a} = 2\sum_{i=1}^{n} (ax_i + b - y_i)x_i = 0\\ \frac{\partial f}{\partial b} = 2\sum_{i=1}^{n} (ax_i + b - y_i) = 0 \end{cases}$$

化简可得解为:

$$\begin{cases} a = \frac{n \sum_{i=1}^{n} x_{i} y_{i} - \sum_{i=1}^{n} x_{i} \sum_{i=1}^{n} y_{i}}{n \sum_{i=1}^{n} x_{i}^{2} - \left(\sum_{i=1}^{n} x_{i}\right)^{2}} \\ b = \frac{1}{n} \sum_{i=1}^{n} y_{i} - \frac{a}{n} \sum_{i=1}^{n} x_{i} \end{cases}$$



#### 最小二乘法的矩阵解法 (列空间角度)

首先记:

$$\mathbf{x} = \begin{bmatrix} x_1 & x_2 & \cdots & x_n \end{bmatrix}^{\mathrm{T}}, \ \mathbf{y} = \begin{bmatrix} y_1 & y_2 & \cdots & y_n \end{bmatrix}^{\mathrm{T}}, \ \mathbf{1} = \begin{bmatrix} 1 & 1 & \cdots & 1 \end{bmatrix}^{\mathrm{T}}$$

则线性模型可以表示为:

$$\mathbf{y} = a\mathbf{x} + b\mathbf{l}$$

进一步, 令

$$\mathbf{X} = \begin{bmatrix} \mathbf{x} & \mathbf{l} \end{bmatrix} \qquad \mathbf{c} = \begin{bmatrix} a & b \end{bmatrix}^{\mathrm{T}}$$

则线性模型可以进一步表示为:

$$y = Xc$$



优化模型转化为:

$$\min_{\mathbf{c}} f(\mathbf{c}) = \left\| \mathbf{y} - \mathbf{X} \mathbf{c} \right\|^2$$

对自变量求导,并令其为0

$$f_{\mathbf{c}}' = 2\mathbf{X}^{\mathrm{T}}\mathbf{X}\mathbf{c} - 2\mathbf{X}^{\mathrm{T}}\mathbf{y} = \mathbf{0}$$

可得模型的解为

$$\mathbf{c} = \left(\mathbf{X}^{\mathrm{T}}\mathbf{X}\right)^{-1}\mathbf{X}^{\mathrm{T}}\mathbf{y} \qquad \mathbf{X}^{\#} = \left(\mathbf{X}^{\mathrm{T}}\mathbf{X}\right)^{-1}\mathbf{X}^{\mathrm{T}}$$



#### 一个简单的例子: 假设有三个观测点, 分别为

(1,1),(2,1)和(3,3)

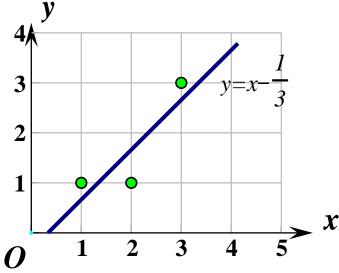
#### 方法1:

$$a = \frac{n\sum_{i=1}^{n} x_{i} y_{i} - \sum_{i=1}^{n} x_{i} \sum_{i=1}^{n} y_{i}}{n\sum_{i=1}^{n} x_{i}^{2} - \left(\sum_{i=1}^{n} x_{i}\right)^{2}} = \frac{3(1*1+2*1+3*3) - (1+2+3)(1+1+3)}{3(1^{2}+2^{2}+3^{2}) - (1+2+3)^{2}} = 1$$

$$b = \frac{1}{n} \sum_{i=1}^{n} y_i - \frac{a}{n} \sum_{i=1}^{n} x_i = \frac{1}{3} (1 + 1 + 3) - \frac{1}{3} (1 + 2 + 3) = -\frac{1}{3}$$

#### 方法2:

$$\mathbf{c} = \left(\mathbf{X}^{\mathsf{T}}\mathbf{X}\right)^{-1}\mathbf{X}^{\mathsf{T}}\mathbf{y} = \begin{bmatrix} 1 & 2 & 3 \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 2 & 1 \\ 3 & 1 \end{bmatrix} \begin{bmatrix} 1 & 2 & 3 \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 3 \end{bmatrix} = \begin{bmatrix} 1 \\ -\frac{1}{3} \end{bmatrix}$$



### 变量问题

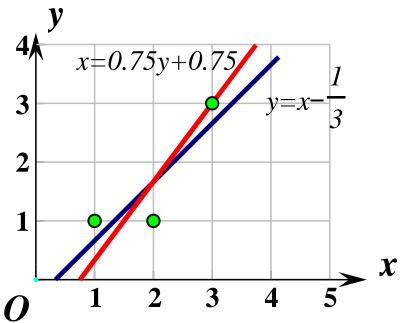


#### 最小二乘法中的变量问题

在前面,我们通过给出一组散点的最佳直线拟合阐述了最小二乘法的基本原理。对于同样的一组散点,我们用x = a'y + b'来拟合这组散点是否可以得到同样的结果呢?

$$\begin{cases} x_1 = a' y_1 + b' \\ x_2 = a' y_2 + b' \\ \vdots \\ x_n = a' y_n + b' \end{cases}$$

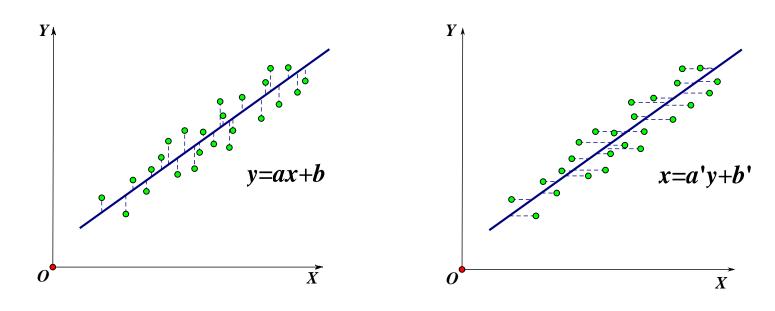
$$\mathbf{c'} = \left(\mathbf{Y}^{\mathrm{T}}\mathbf{Y}\right)^{-1}\mathbf{Y}^{\mathrm{T}}\mathbf{x} = \begin{bmatrix} 0.75\\0.75 \end{bmatrix}$$



### 变量问题



#### 最小二乘法中的变量问题

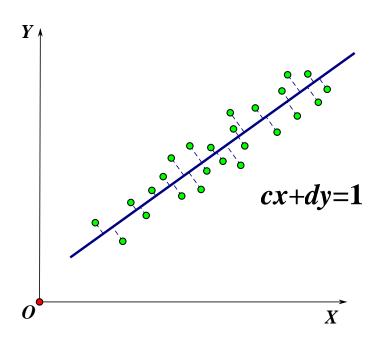


当x和y互为自变量和因变量直线拟合的直观解释

### 变量问题



#### 总体最小二乘法

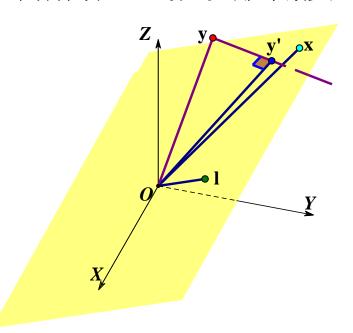


总体最小二乘法的直观理解

### 几何解释



#### 最小二乘法的几何解释(线性变换角度)



最小二乘法对应的样本空间的几何: 其中 $\mathbf{x}$ =[1 2 3]<sup>T</sup>,  $\mathbf{y}$ =[1 1 3]<sup>T</sup>,  $\mathbf{1}$ =[1 1 1]<sup>T</sup>。最小二乘揭示的是三个散点之间的关系。当以 $\mathbf{y}$ = $\mathbf{a}\mathbf{x}$ + $\mathbf{b}$  来拟合散点时,对应到样本空间,相当于寻求向量 $\mathbf{y}$ 到 $\mathbf{x}$ 和 $\mathbf{1}$ 这两个向量所张成平面的投影

### 几何解释



#### 最小二乘法的几何解释(线性变换角度)

$$\min_{\mathbf{c}} f(\mathbf{c}) = \left\| \mathbf{y} - \mathbf{X} \mathbf{c} \right\|^2$$

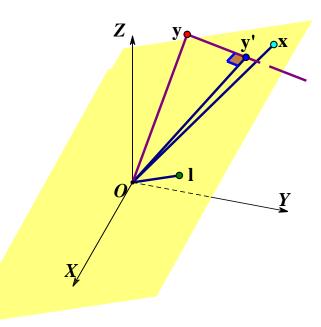
$$\mathbf{c} = \left(\mathbf{X}^{\mathrm{T}}\mathbf{X}\right)^{-1}\mathbf{X}^{\mathrm{T}}\mathbf{y}$$

$$\mathbf{y'} = \mathbf{X}\mathbf{c} = \mathbf{X}\left(\mathbf{X}^{\mathrm{T}}\mathbf{X}\right)^{-1}\mathbf{X}^{\mathrm{T}}\mathbf{y}$$

矩阵的广义逆:  $\mathbf{X}^{\#} = \left(\mathbf{X}^{\mathrm{T}}\mathbf{X}\right)^{-1}\mathbf{X}^{\mathrm{T}}$ 

投影矩阵 
$$\mathbf{P}_{\mathbf{X}} = \mathbf{X} \left( \mathbf{X}^{\mathrm{T}} \mathbf{X} \right)^{-1} \mathbf{X}^{\mathrm{T}} = \mathbf{X} \mathbf{X}^{\#}$$
  $\mathbf{y}' = \mathbf{P}_{\mathbf{x}} \mathbf{y}$ 

正交补投影算子 
$$\mathbf{P}_{\mathbf{x}}^{\perp} = \mathbf{I} - \mathbf{P}_{\mathbf{x}} = \mathbf{I} - \mathbf{X}\mathbf{X}^{\#}$$

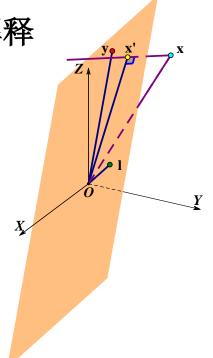


讨论: 投影矩阵的性质?

### 几何解释



最小二乘法的几何解释



最小二乘法对应的样本空间的几何:其中 $\mathbf{x}=[1\ 2\ 3]^T$ ,  $\mathbf{y}=[1\ 1\ 3]^T$ ,  $\mathbf{1}=[1\ 1\ 1]^T$ 。最小二乘揭示的是三个散点之间的关系。当以 $\mathbf{x}=a'y+b'$ 来拟合散点时,对应到样本空间,相当于寻求 $\mathbf{x}$ 到 $\mathbf{y}$ 和 $\mathbf{1}$ 这两个向量所张成平面的投影



#### 最小二乘法的概率解释

在前面的讨论中, 因变量的总体观测误差用

$$f(a,b) = \sum_{i=1}^{n} (ax_i + b - y_i)^2$$

来表示。即对于每个观测点的因变量的观测误差,都选用模型解与观测值的差的平方来衡量。这正是相应的方法命名为最小二乘法而不是最小一乘或者最小三乘的原因所在。

### Why?



#### 最小二乘法的概率解释

记

$$\begin{cases} y_1 = \mathbf{c}^T \mathbf{x}_1 + \varepsilon_1 \\ y_2 = \mathbf{c}^T \mathbf{x}_2 + \varepsilon_2 \\ \vdots \\ y_n = \mathbf{c}^T \mathbf{x}_n + \varepsilon_n \end{cases}$$

$$\mathbf{c} = \begin{bmatrix} a & b \end{bmatrix}^T$$

$$\mathbf{x}_i = \begin{bmatrix} x_i & 1 \end{bmatrix}^T$$

其中 $\varepsilon_i$ 对应第i个因变量 $y_i$ 的观测误差,它对应着不能被线性模型刻画的因素。假设 $\varepsilon_i$ 服从正态分布,

$$\varepsilon_i \sim N(0, \sigma^2)$$

则其概率密度函数为,

$$f\left(\varepsilon_{i}\right) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{\varepsilon_{i}^{2}}{2\sigma^{2}}\right)$$



这意味着,在给定 $\mathbf{x}_i$ 和参数 $\mathbf{c}$ 的情况下,因变量 $y_i$ 也服从正态分布,即

$$f(y_i|\mathbf{x}_i;\mathbf{c}) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y_i - \mathbf{c}^{\mathrm{T}}\mathbf{x}_i)^2}{2\sigma^2}\right) \qquad \mathbf{c} = \begin{bmatrix} a & b \end{bmatrix}^{\mathrm{T}} \\ \mathbf{x}_i = \begin{bmatrix} x_i & 1 \end{bmatrix}^{\mathrm{T}}$$

假设所有的 $\varepsilon_i$ 独立同分布,我们可以定义所有观测数据关于参数 $\mathbf{c}$ 的似然函数如下:

$$L(\mathbf{c}) = f\left(\mathbf{y} \middle| \mathbf{X}; \mathbf{c}\right) = \prod_{i=1}^{n} f\left(y_{i} \middle| \mathbf{x}_{i}; \mathbf{c}\right) = \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{\left(y_{i} - \mathbf{c}^{\mathrm{T}} \mathbf{x}_{i}\right)^{2}}{2\sigma^{2}}\right)$$

其中

$$\mathbf{y} = \begin{bmatrix} y_1 & y_2 & \cdots & y_n \end{bmatrix}^{\mathrm{T}}, \ \mathbf{X} = \begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \cdots & \mathbf{x}_n \end{bmatrix}^{\mathrm{T}}$$



为了便于求解,定义  $l(\mathbf{c}) = \ln(L(\mathbf{c}))$ 

$$\begin{split} &l(\mathbf{c}) = \ln\left(L\left(\mathbf{c}\right)\right) = \ln\left(\prod_{i=1}^{n} f\left(y_{i} \middle| \mathbf{x}_{i}; \mathbf{c}\right)\right) \\ &= \ln\left(\prod_{i=1}^{n} \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{\left(y_{i} - \mathbf{c}^{\mathrm{T}} \mathbf{x}_{i}\right)^{2}}{2\sigma^{2}}\right)\right) \\ &= \sum_{i=1}^{n} \ln\left(\frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{\left(y_{i} - \mathbf{c}^{\mathrm{T}} \mathbf{x}_{i}\right)^{2}}{2\sigma^{2}}\right)\right) \\ &= n \ln\left(\frac{1}{\sqrt{2\pi\sigma}}\right) - \sum_{i=1}^{n} \frac{\left(y_{i} - \mathbf{c}^{\mathrm{T}} \mathbf{x}_{i}\right)^{2}}{2\sigma^{2}} \end{split}$$

因此,最大化似然函数,相当于最小化

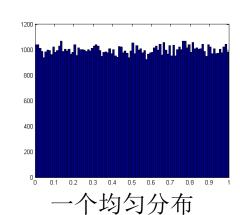
$$\sum_{i=1}^{n} \left( y_i - \mathbf{c}^{\mathrm{T}} \mathbf{x}_i \right)^2 = \sum_{i=1}^{n} \left( a x_i + b - y_i \right)^2 = f(a, b)$$

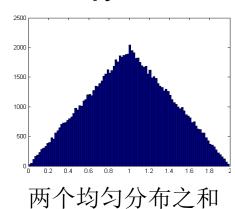


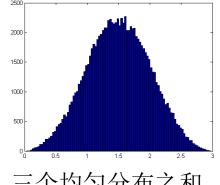
从上面的推导可以看出,各个观测点的模型误差满足独立 同分布的高斯分布是最小二乘法能够行之有效的概率机制 前提。如果此条件不能得到满足,用最小二乘法拟合数据 将不能得到最优解。而根据**中心极限定理**,这个条件一般 情况下是近似成立的。

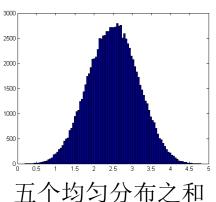
中心极限定理:对于多个独立的随机变量,它们和的平均,满足

$$\lim_{n\to\infty}\frac{x_1+x_2+\cdots+x_n}{n}\to \text{ass}$$









三个均匀分布之和

中国科学院空间信息处理与应用系统技术重点实验室

### 约束最小二乘



▶ 讨论: 四种模型所对应的几何投影

$$\min_{\mathbf{c}} \|\mathbf{y} - \mathbf{X}\mathbf{c}\|^2$$
 无约束

$$\begin{cases} \min_{\mathbf{c}} \|\mathbf{y} - \mathbf{X}\mathbf{c}\|^2 \\ \mathbf{c}^{\mathbf{T}}\mathbf{1} = 1 \end{cases}$$
 \(\frac{\pm \text{\$\frac{1}{2}\$}}{\pm \text{\$\frac{1}{2}\$}}

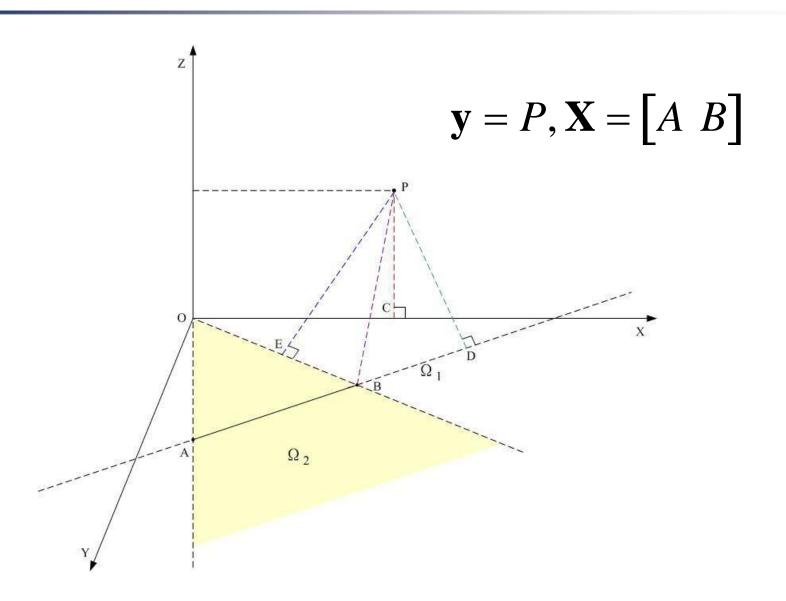
$$\begin{cases} \min_{\mathbf{c}} \|\mathbf{y} - \mathbf{X}\mathbf{c}\|^2 \\ \mathbf{c} \ge \mathbf{0} \end{cases}$$

$$\begin{cases} \min_{\mathbf{c}} \|\mathbf{y} - \mathbf{X}\mathbf{c}\|^2 \\ \mathbf{c}^{\mathrm{T}}\mathbf{1} = 1 \end{cases} \stackrel{\text{$\pm 0$}}{}$$

其中1为所有元素都为1的列向量,0为所有元素都为0的列向量

## 约束最小二乘



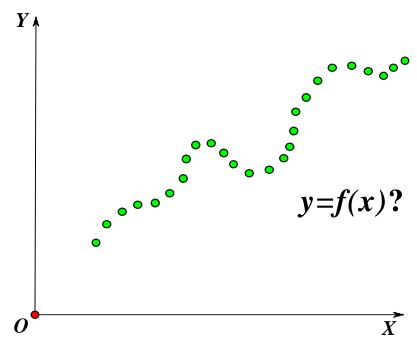




#### 讨论:

➤ 除了可以拟合直线,最小二乘可以拟合任意次数的多项式,乃至任意已知表达式的曲线。那么如何用最小二乘法拟合任意形状的未知表达式的曲线呢?

$$y = f(x) = x^{n} + a_{1}x^{n-1} + \dots + a_{n-1}x + a_{n}$$



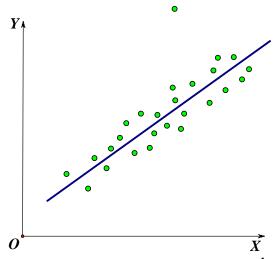


#### 思考题(每题1分,共3分):

▶ 病态情形: 当X<sup>T</sup>X 不可逆时,如何处理?

$$\mathbf{c} = \left(\mathbf{X}^{\mathsf{T}} \mathbf{X}\right)^{-1} \mathbf{X}^{\mathsf{T}} \mathbf{y}$$

▶ 异常情形: 当数据中存在异常点时,如何处理?



中国科学院空间信息处理与应用系统技术重点实验室



#### 思考题(每题1分,共3分):

> 求解如下等式约束的最小二乘解

$$\begin{cases} \min_{\mathbf{c}} \|\mathbf{y} - \mathbf{X}\mathbf{c}\|^2 \\ s.t. \ \mathbf{c}^{\mathsf{T}}\mathbf{1} = 1 \end{cases}$$

其中为1所有元素都为1的列向量



#### 最小二乘法的遥感应用:

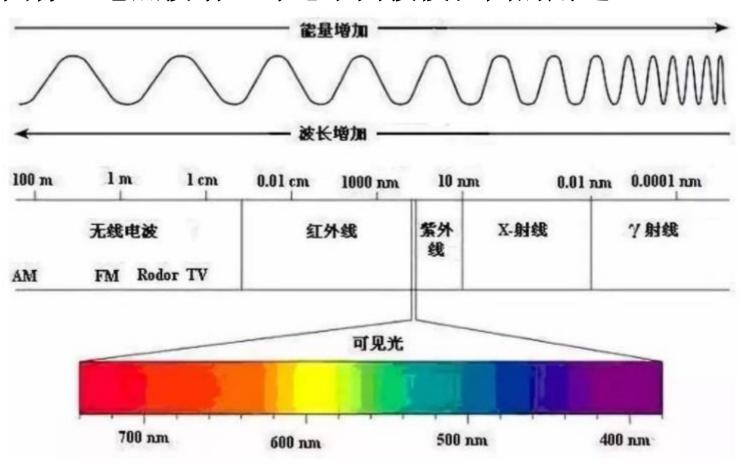
遥感图像简介

遥感图像混合像元分析

遥感图像条带消除

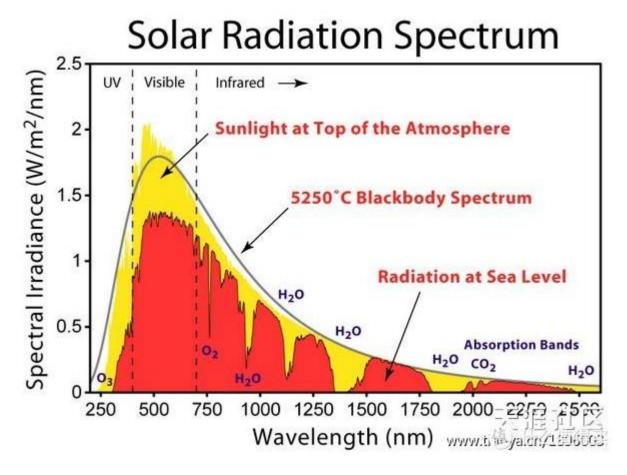


遥感图像、电磁波谱(讨论不同波段范围的用途?)





#### 太阳辐射光谱





#### 遥感图像发展趋势:

空间分辨率

光谱分辨率

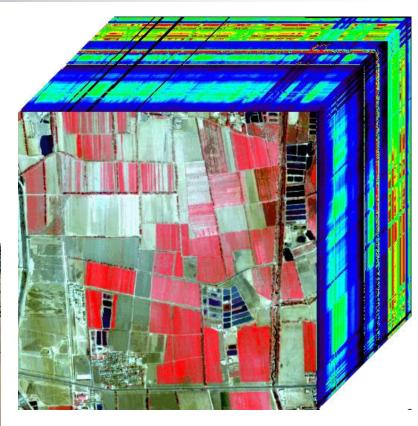
时间分辨率

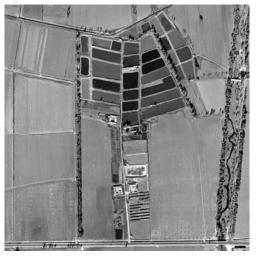
波谱范围



全色-(彩色)-多光谱-高光谱

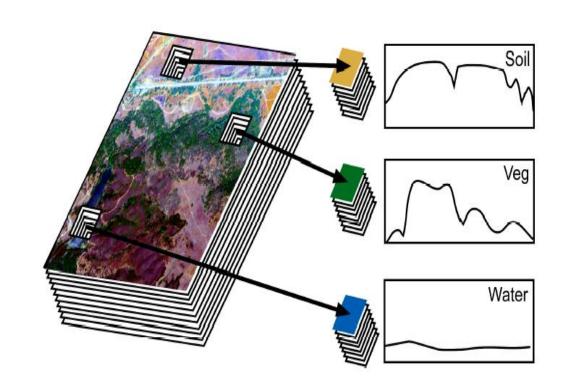








### 高光谱遥感的概念(ENVI演示)

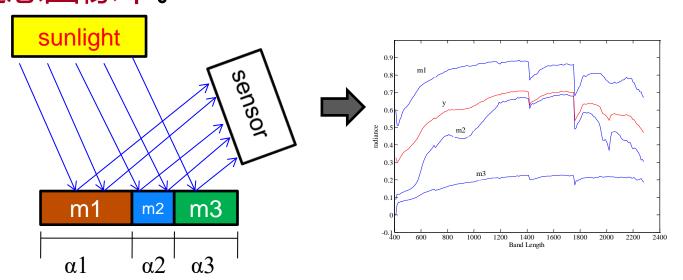


从每个象元均可提取一条连续的光谱曲线



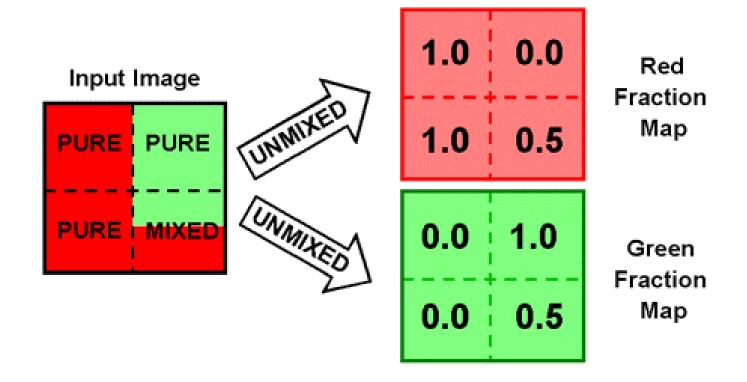
### 混合像元分析

遥感图象中每个象元所对应的地表,往往包含不同的覆盖类型,他们有着不同的光谱响应特征。由于传感器分辨率的限制,混合像元普遍存在于遥感图像中。





### 混合像元分析





### 混合像元分析

图象上的线性混合模型

光谱空间的单形体结构



□线性混合模型(灰度图像情形)

$$p = \sum_{i=1}^{N} c_i e_i + n$$

- 一个方程, N个未知量, 方程有无穷解, 但解都没意义
- □线性混合模型(多光谱图像情形,波段数M)

$$\mathbf{p} = \sum_{i=1}^{N} c_i \mathbf{e}_i + \mathbf{n} = \mathbf{E}\mathbf{c} + \mathbf{n}$$

M个方程,N个未知量,方程组解的情况取决于M和N的大小

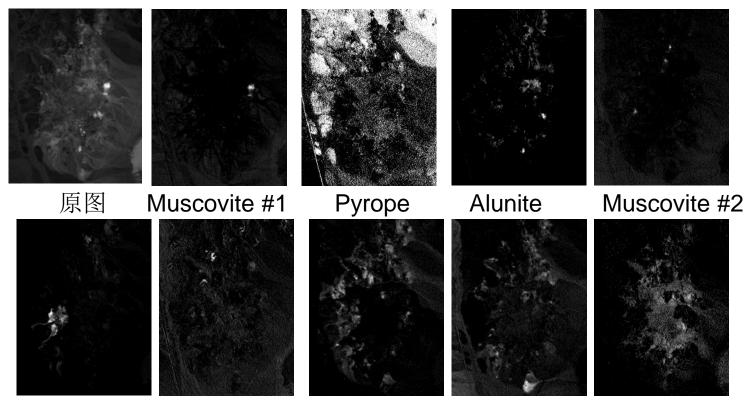


□线性混合模型(高光谱图像情形,波段数*M*)

$$\mathbf{p} = \sum_{i=1}^{N} c_i \mathbf{e}_i + \mathbf{n} = \mathbf{E}\mathbf{c} + \mathbf{n}$$

对于高光谱图像,一般M远大于N,方程有唯一的最小二乘解



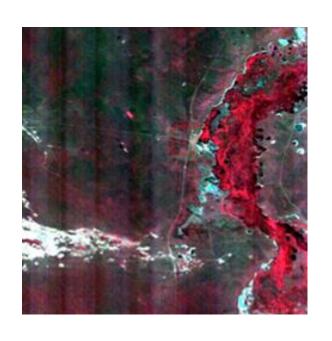


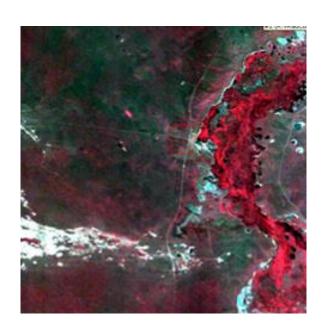
Buddingtonite Nontronite #1 Dumortierite Kaolinite Nontronite #2

美国cuprite区域Aviris高光谱数据混合像元分解部分结果图



### 条带消除

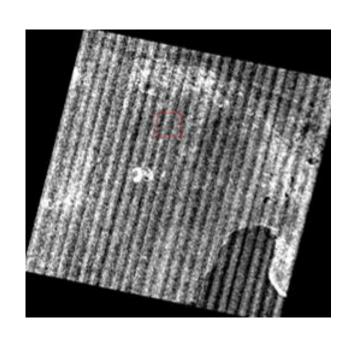


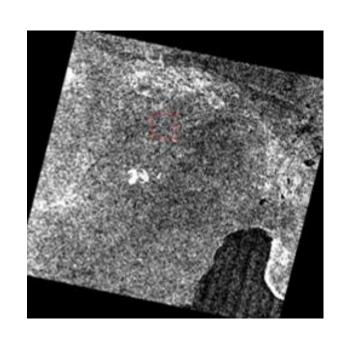


VRadCor(赵永超)批处理HJ-1A已经地面处理系统相对辐射校正后的高光谱单景图像中残存条纹的前(左)后(右)效果对比



### 条带消除





VRadCor(赵永超)批处理HJ-1A已经几何纠正的高光谱单景图像时的前(左)后(右)效果对比,只显示了其中一个波段



### 条带消除





南佛罗里达大学和NOAA在墨西哥湾的海草遥感监测决策系统中采用VRadCor对有限海域HYPERION图像中残存的高频条纹进行修正的效果(赵永超)



### 条带消除

快舟1号利用 VRadCor计算的 定标系数处理前 (上)后(下)的陆域 效果,消除了其 中的低频条带条 纹(赵永超)





中国科学院空间信息处理与应用系统技术重点实验室

### 总结



- 1. 从行空间角度,线性方程组的解为各个方程所对应的超平面的交集
- 2. 从列空间角度,线性方程组的常数项向量为系数矩阵的各个列向量的线性表出,且表出系数即为待求的解
- 3. 最小二乘法是求解矛盾方程组的常用方法
- 4. 代数上,线性方程组的最小二乘解对应着矩阵的广义逆操作
- 5. 几何上,线性方程组的最小二乘解等价于线性方程组的常数项向量 在系数矩阵的列空间的正交投影
- 6. 概率上,最小二乘法基于线性方程组模型误差的高斯分布。



# 谢谢

## 耿修瑞

中国科学院空天信息创新研究院

gengxr@sina.com.cn