

北京市大数据人才培训示范基地

## 第二讲：回归



欧高炎，北京大学博士、博士后  
数据酷客创始人，博雅大数据学院院长

时间	主题	介绍
5/14	机器学习介绍	机器会学习吗？
5/21	回归	回归初心，方得始终
5/28	分类	分门别类，各得其所
6/4	模型提升	三个臭皮匠，顶个诸葛亮
6/11	聚类	物以类聚，人以群分
6/18	降维	取其精华，去其糟粕
6/23	最优化	摸着石头过河，蒙着眼睛爬山
7/2	维度灾难	来自维数的诅咒
7/9	深度学习	深层次学习的艺术
7/16	强化学习	让机器像人类一样学习



- 对于  $n \times n$  方阵  $\mathbf{A}$ ，如果存在矩阵  $\mathbf{B}$  使得  $\mathbf{AB} = \mathbf{BA} = \mathbf{I}$ ，则称  $\mathbf{B}$  为  $\mathbf{A}$  的逆矩阵，记为  $\mathbf{A}^{-1}$
- 若  $\mathbf{A}$  为可逆矩阵，则其逆矩阵是唯一的
- 如何判断矩阵是否可逆？
  - 行列式不等于0
  - 满秩
  - 行（或列）向量组线性无关
  - ...

- numpy.linalg 模块包含线性代数的函数，可计算逆矩阵、求特征值、解线性方程组以及求解行列式等。

- 行列式：np.linalg.det(A)

- 计算逆矩阵：np.linalg.inv(A)

```
import numpy as np
#格式化numpy输出
np.set_printoptions(formatter={'float': '{: 0.2f}'.format})
A = np.array([[4,5,1],[0,8,3],[9,4,7]])
#求行列式
print("行列式: " + str(np.linalg.det(A)))
B = np.linalg.inv(A) #求逆
print("A的逆矩阵为B: ")
print(A_inv)
print("BA:") #结果验证
print(B.dot(A))
print("AB:")
print(A.dot(B))
```

行列式: 239.00000000000009

A的逆矩阵为B:

```
[[ 0.18 -0.13  0.03]
 [ 0.11  0.08 -0.05]
 [-0.30  0.12  0.13]]
```

BA:

```
[[ 1.00  0.00  0.00]
 [ 0.00  1.00  0.00]
 [-0.00  0.00  1.00]]
```

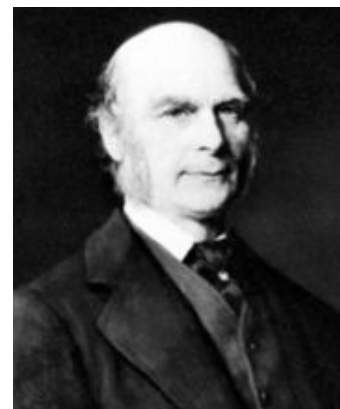
AB:

```
[[ 1.00 -0.00  0.00]
 [ 0.00  1.00  0.00]
 [ 0.00  0.00  1.00]]
```

- 回归 ( Regression ) 这一概念最早由英国生物统计学家高尔顿和他的学生皮尔逊在研究父母亲和子女的身高遗传特性时提出
- “子女的身高趋向于高于父母的身高的平均值，但一般不会超过父母的身高。”-- 《遗传的身高向平均数方向的回归》

$$Y = 0.8567 + 0.516 * X \quad (\text{单位为米})$$

- 回归如今指的用一个或多个自变量来预测因变量的数学方法
- 在机器学习中，回归指的是一类预测变量为连续值的有监督学习方法



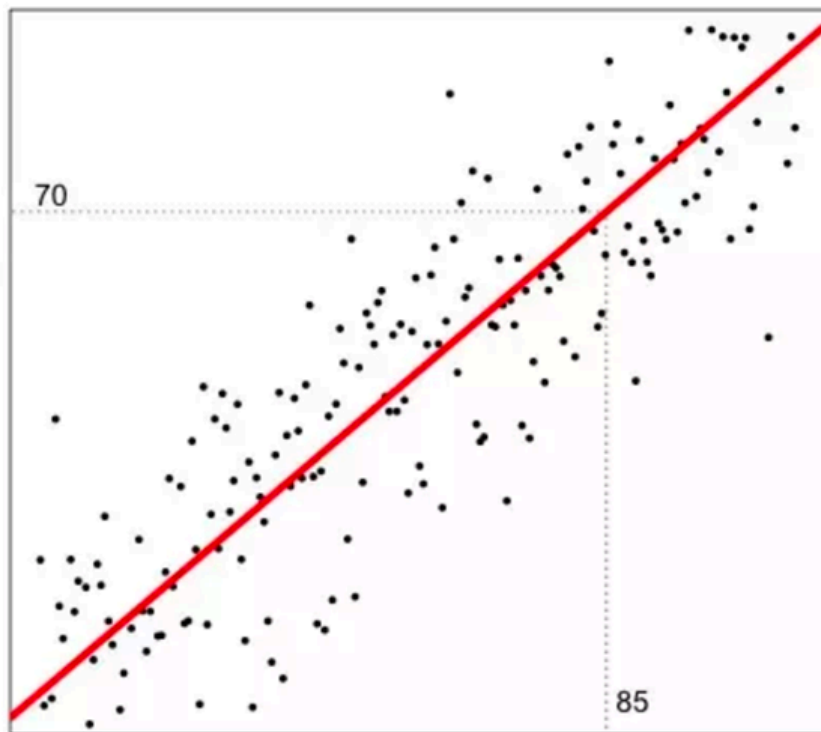
高尔顿



皮尔逊

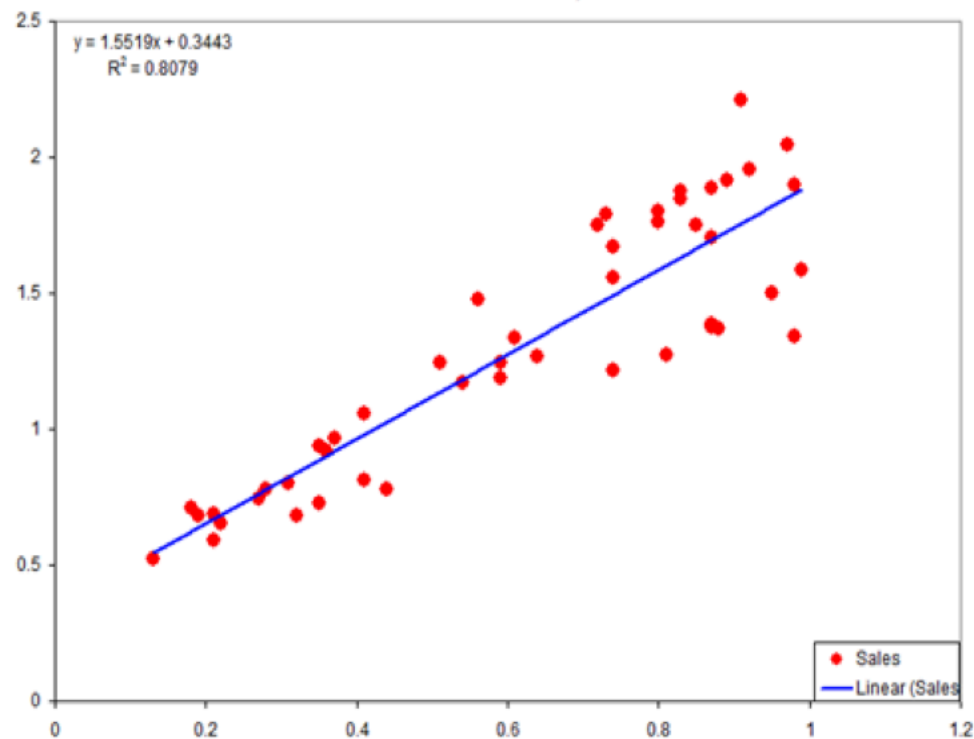
在回归模型中，需要预测的变量叫做**因变量**，用来解释因变量变化的变量叫做**自变量**。

成为  
领导者可  
能性



IQ

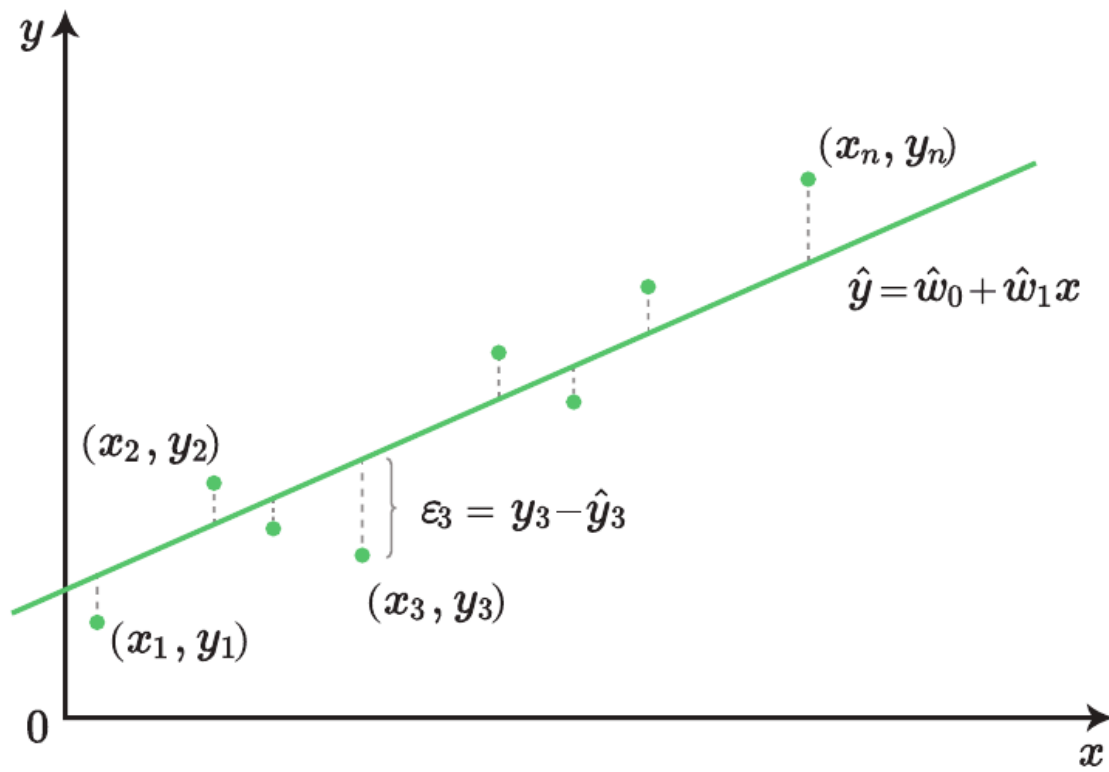
销量



广告投入

- 模型为  $y = w_1x + w_0$  , 其中  $w_0, w_1$  为回归系数
- 给定训练集  $D = \{(x_1, y_1), \dots, (x_n, y_n)\}$  , 我们的目标是找到一条直线  $y = w_1x + w_0$  使得所有样本尽可能落在它的附近。
- 优化目标为：

$$\min_{(w_0, w_1)} \sum_{i=1}^n (y_i - w_1x_i - w_0)^2$$





- 记优化目标为  $L(w_1, w_0) = \sum_{i=1}^n (y_i - w_1 x_i - w_0)^2$
- $L(w_1, w_0)$  为二次凸函数，分别对  $w_1, w_0$  求导并令导数为零：

$$\begin{cases} \frac{\partial L(w_1, w_0)}{\partial w_1} = 2 \sum_{i=1}^n (y_i - w_1 x_i - w_0)(-x_i) = 0 \\ \frac{\partial L(w_1, w_0)}{\partial w_0} = 2 \sum_{i=1}^n (y_i - w_1 x_i - w_0)(-1) = 0 \end{cases}$$

- 解上述线性方程组得

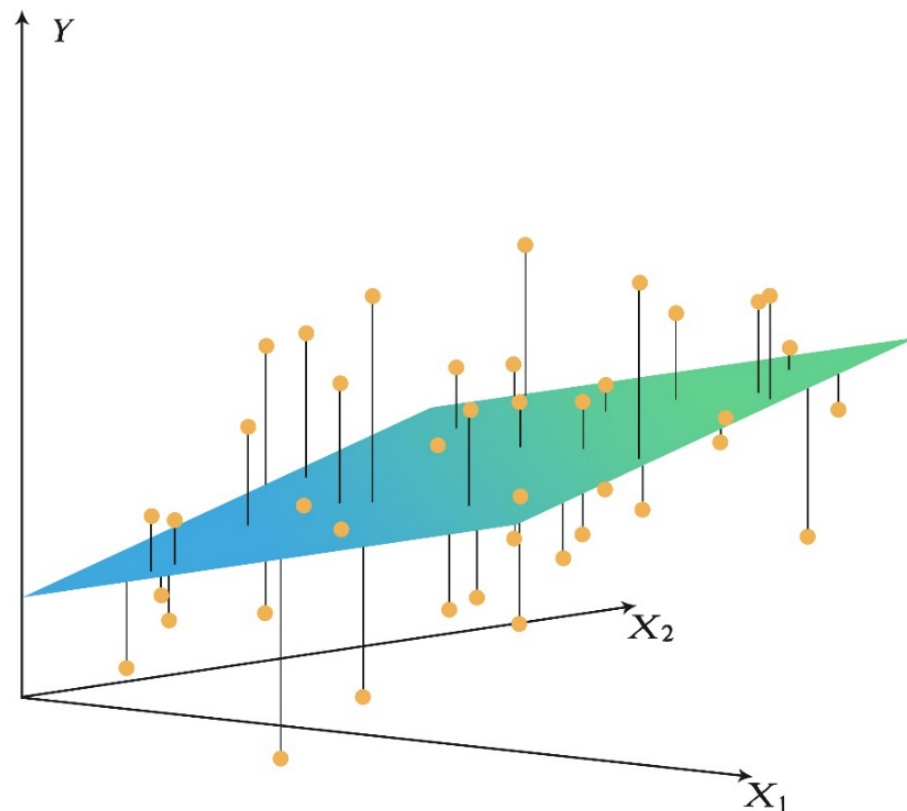
$$\begin{cases} w_1 = \frac{n \sum_{i=1}^n y_i x_i - (\sum_{i=1}^n y_i)(\sum_{i=1}^n x_i)}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2} \\ w_0 = \frac{\sum_{i=1}^n y_i}{n} - \frac{\sum_{i=1}^n x_i}{n} w_1 \end{cases}$$



- $y$  是多个特征的线性组合

$$y = w_1x_1 + w_2x_2 + \cdots + w_dx_d + w_0$$

- 训练集  $D = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$
- $\mathbf{x}_i$  为  $d$  维特征向量  $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{id})^T$
- 假设模型对  $\mathbf{x}_i$  的预测值为  $\hat{y}_i = w_1x_{i1} + \cdots + w_dx_{id} + w_0$
- 优化目标为  $L(\mathbf{w}) = \sum_{i=1}^n (\hat{y}_i - y_i)^2$



寻找一个超平面，使得训练集中样本到超平面的误差平方和最小

- 假设训练集的特征部分记为  $n \times (d + 1)$  矩阵  $\mathbf{X}$  , 其中最后一列取值全为 1
- 标签部分记为  $\mathbf{y} = (y_1, y_2, \dots, y_n)^T$  , 参数记为  $\mathbf{w} = (w_1, w_2, \dots, w_d, w_0)^T$
- 多元线性回归模型为 :

$$\hat{\mathbf{y}} = \mathbf{X}\mathbf{w}$$

$$\begin{bmatrix} x_{11} & x_{12} & \dots & x_{1d} & 1 \\ x_{21} & x_{22} & \dots & x_{2d} & 1 \\ \dots & \dots & \dots & \dots & 1 \\ x_{i1} & x_{i2} & \dots & x_{id} & 1 \\ \dots & \dots & \dots & \dots & \dots \\ x_{n1} & x_{n2} & \dots & x_{nd} & 1 \end{bmatrix}$$

$\mathbf{X}$

$$\begin{bmatrix} w_1 \\ w_2 \\ \dots \\ w_d \\ w_0 \end{bmatrix}$$

$\mathbf{w}$

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \dots \\ \hat{y}_i \\ \dots \\ \hat{y}_n \end{bmatrix}$$

$\hat{\mathbf{y}}$

$$\begin{bmatrix} y_1 \\ y_2 \\ \dots \\ y_i \\ \dots \\ y_n \end{bmatrix}$$

$\mathbf{y}$

- 最小化均方误差函数为： $L(\mathbf{w}) = \sum_{i=1}^n (\hat{y}_i - y_i)^2$

- $$L(\mathbf{w}) = (\hat{\mathbf{y}} - \mathbf{y})^T (\hat{\mathbf{y}} - \mathbf{y})$$
$$= (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y})$$

$$\hat{\mathbf{y}} - \mathbf{y} = \begin{bmatrix} \hat{y}_1 - y_1 \\ \hat{y}_2 - y_2 \\ \dots \\ \hat{y}_i - y_i \\ \dots \\ \hat{y}_n - y_n \end{bmatrix}$$

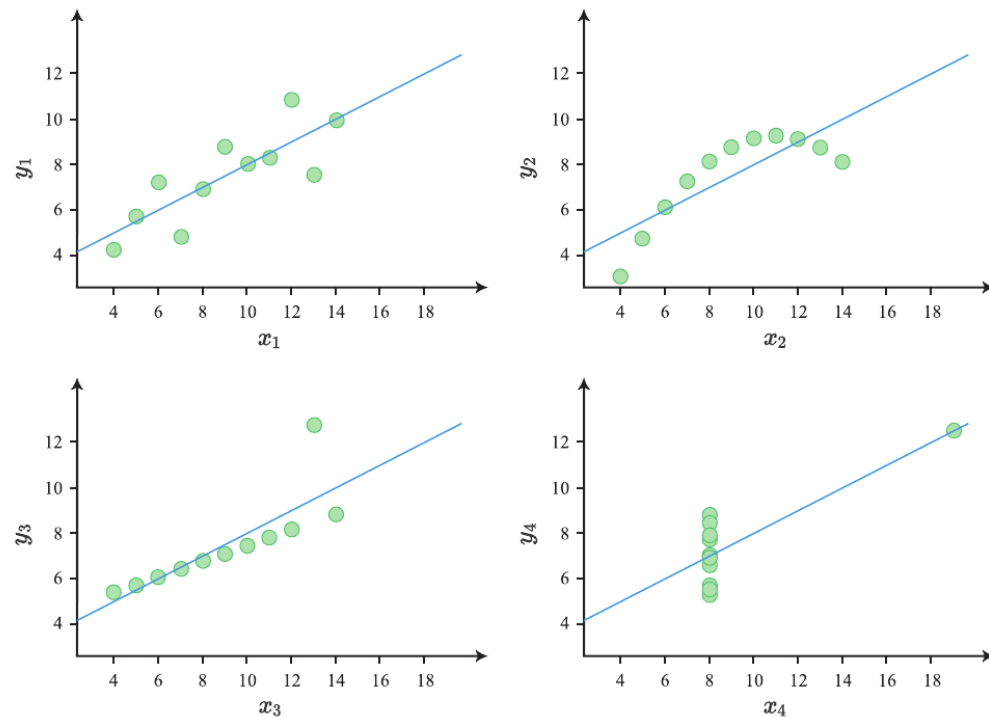
$$(\hat{\mathbf{y}} - \mathbf{y})^T (\hat{\mathbf{y}} - \mathbf{y}) = \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

- 优化目标函数为： $L(\mathbf{w}) = (\mathbf{X}\mathbf{w} - \mathbf{y})^T(\mathbf{X}\mathbf{w} - \mathbf{y})$
- 当矩阵  $\mathbf{X}^T\mathbf{X}$  满秩时：

$$\text{令 } \frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} = -2\mathbf{X}^T(\mathbf{y} - \mathbf{X}\mathbf{w}) = \mathbf{0}, \text{ 可得： } \hat{\mathbf{w}} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$$

## 安斯库姆四重奏

- 实际数据可能不是线性的？
- 多项式回归：使用原始特征的二次项、三次项
  - 线性回归解决非线性问题
  - 问题：维度灾难、过度拟合



- 多重共线性

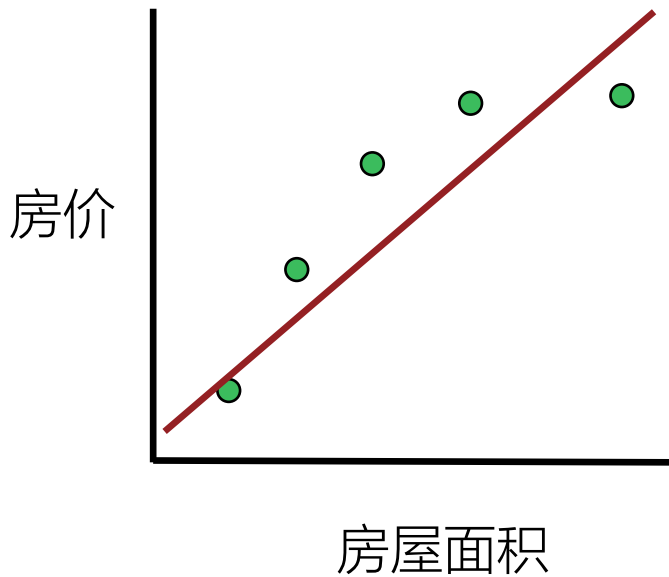
最小二乘的参数估计为  $\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ ，如果变量之间存在较强的共线性，则  $\mathbf{X}^T \mathbf{X}$  近似奇异，对参数的估计变得不准确，造成过度拟合现象。

- 解决方法：正则化、主成分回归、偏最小二乘回归

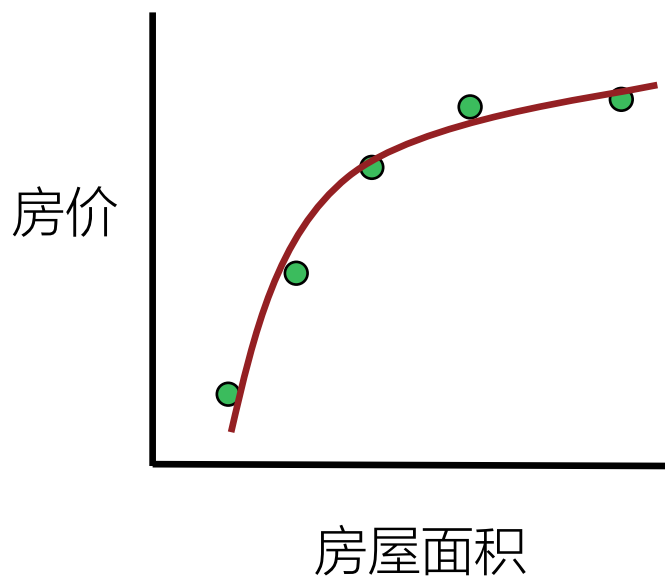
- 过度拟合问题：当模型的变量过多时，线性回归可能会出现过度拟合问题

例如在房价预测问题中，假设 $x$ 表示房屋面积，如果将 $x^2, x^3$ 等作为变量引入可能出现如下情况：

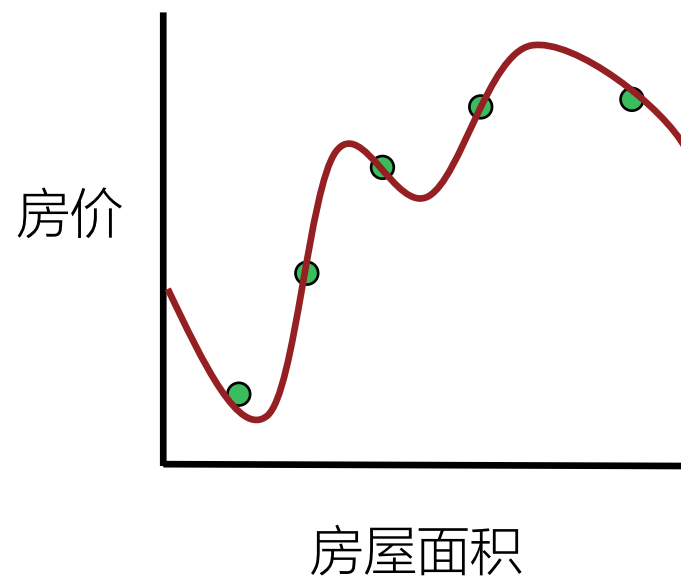
$w_0 + w_1x$ ：  
训练误差大



$w_0 + w_1x + w_2x^2$ ：  
拟合很好



$w_0 + w_1x + w_2x^2 + w_3x^3 + w_4x^4$ ：  
过度拟合



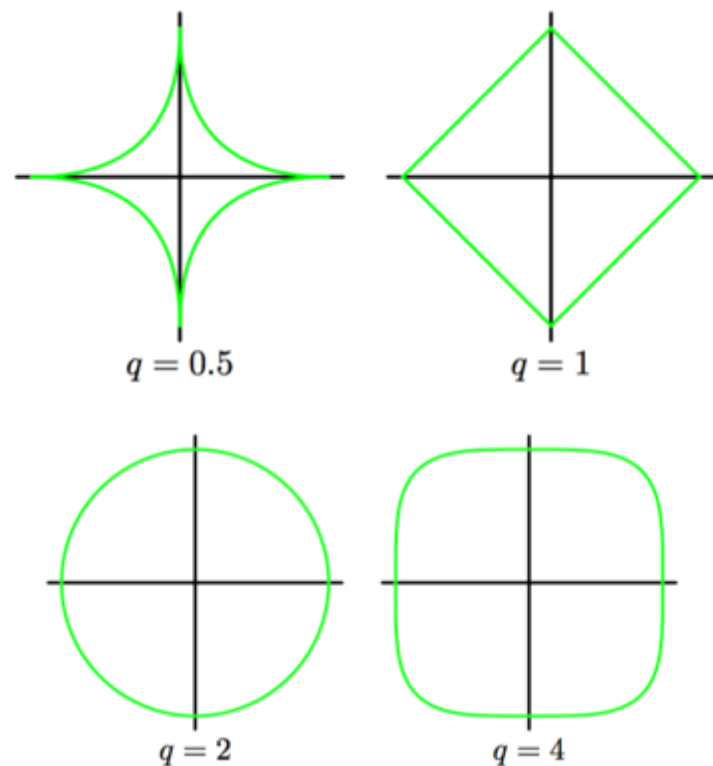


- 正则化可以减小线性回归的过度拟合和多重共线性等问题

$$(\mathbf{X}\mathbf{w} - \mathbf{y})^T(\mathbf{X}\mathbf{w} - \mathbf{y}) + \lambda \|\mathbf{w}\|_q^q$$

正则项或惩罚函数

- $q = 2$  : 岭回归 ( Ridge )
- $q = 1$  : LASSO



$\|\mathbf{w}\|_q$  示意图

- 岭回归思路：在最小二乘法的目标函数上加上一个对 $\mathbf{w}$ 的惩罚函数

$$\lambda \|\mathbf{w}\|_2^2$$

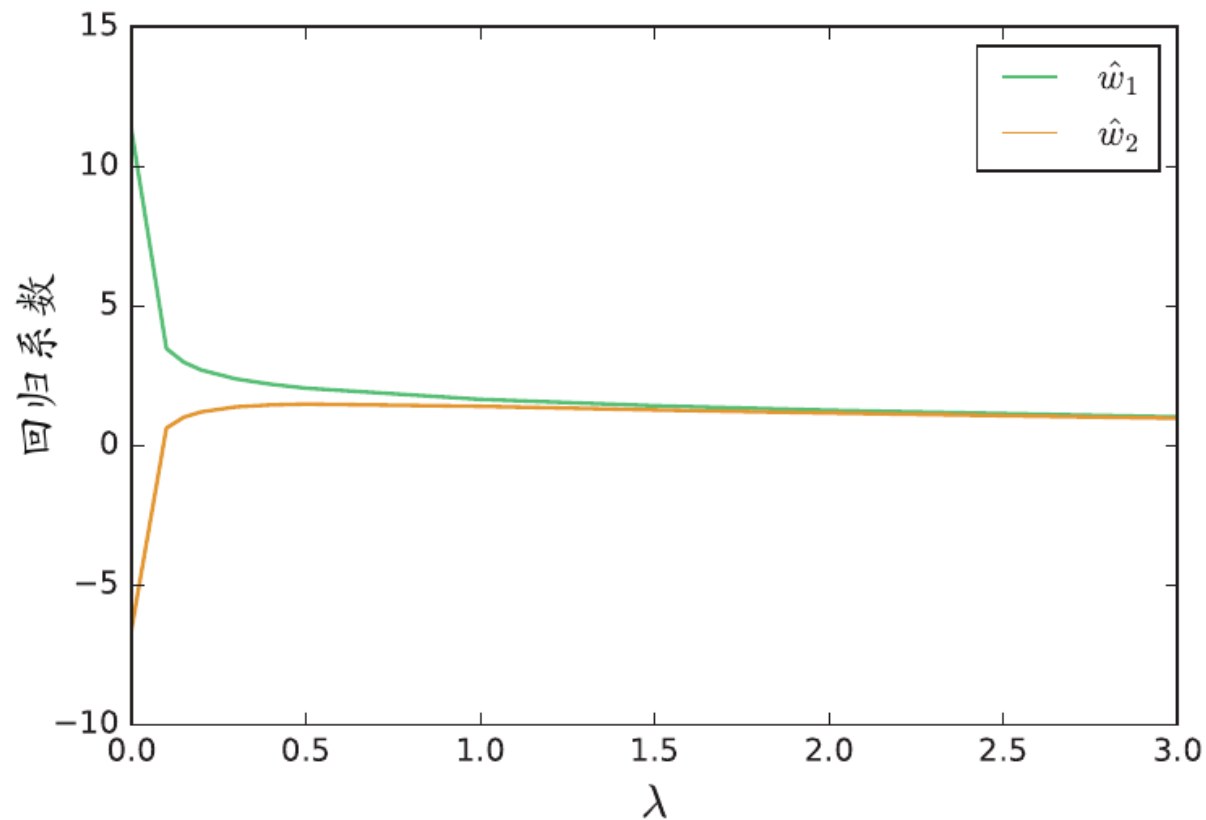
- 线性回归的目标函数： $(\mathbf{X}\mathbf{w} - \mathbf{y})^T(\mathbf{X}\mathbf{w} - \mathbf{y})$
- 岭回归的目标函数变为： $(\mathbf{X}\mathbf{w} - \mathbf{y})^T(\mathbf{X}\mathbf{w} - \mathbf{y}) + \lambda \|\mathbf{w}\|_2^2$
- 对  $\mathbf{w}$  求导并令导数等于零易得： $\hat{\mathbf{w}}^{\text{ridge}} = (\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I})^{-1}\mathbf{X}^T\mathbf{y}$

- 当不断增大正则化参数 $\lambda$ ，估计参数 $\hat{\mathbf{w}}^{\text{ridge}}(\lambda)$ （也称岭回归系数）在坐标系上的变化曲线称为岭迹。岭迹波动很大，说明该变量有共线性。

- 下表是一个二元岭回归的结果：

$\lambda$	0	0.1	0.15	0.2	0.3	0.4	0.5	1.0	1.5	2.0	3.0
$\hat{\mathbf{w}}_1^{\text{ridge}}(\lambda)$	11.31	3.48	2.99	2.71	2.39	2.20	2.06	1.66	1.43	1.27	1.03
$\hat{\mathbf{w}}_2^{\text{ridge}}(\lambda)$	-6.59	0.63	1.02	1.21	1.39	1.46	1.49	1.41	1.28	1.17	0.98

- 根据岭迹做超参数  $\lambda$  的选择
- 在  $\lambda \in (0, 0.5)$  的范围内波动较大，故需要加入正则化项重新进行参数估计，可选  $\lambda = 1$

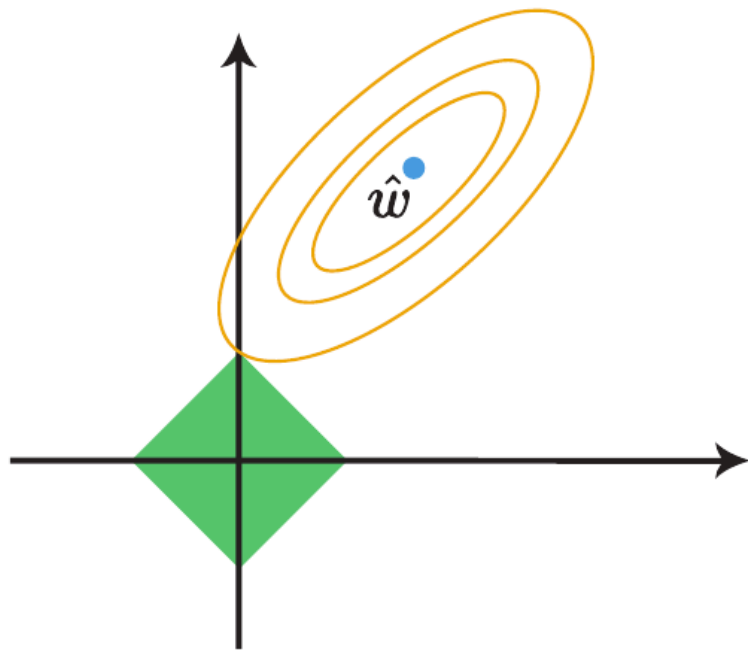


- LASSO 是一种系数压缩估计方法，它的基本思想是通过追求稀疏性自动选择重要的变量
- LASSO 的目标函数： $(\mathbf{X}\mathbf{w} - \mathbf{y})^T(\mathbf{X}\mathbf{w} - \mathbf{y}) + \lambda\|\mathbf{w}\|_1$
- LASSO 的解  $\hat{\mathbf{w}}^{\text{LASSO}}$  没有解析表达式，常用的求解算法包括坐标下降法、LARS算法和 ISTA算法等

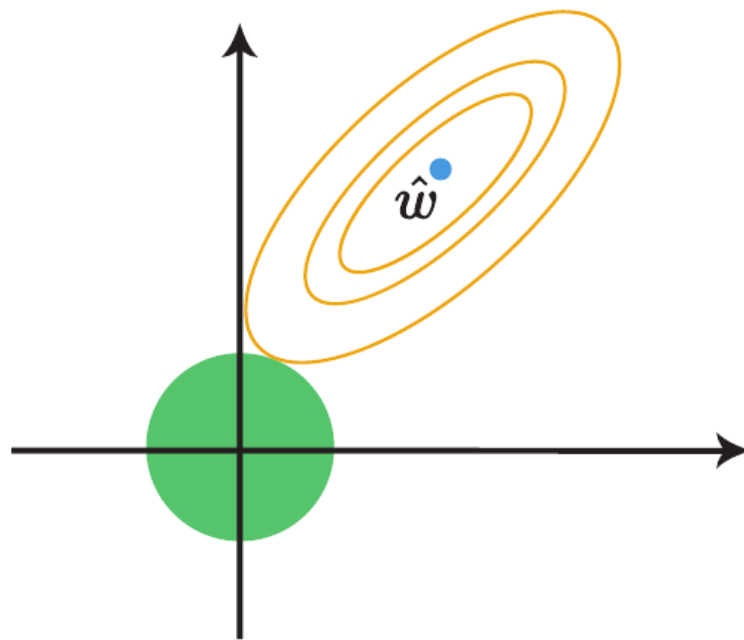
岭回归与 LASSO 如何寻求最优解？

惩罚函数（绿色）与目标函数（橙色）

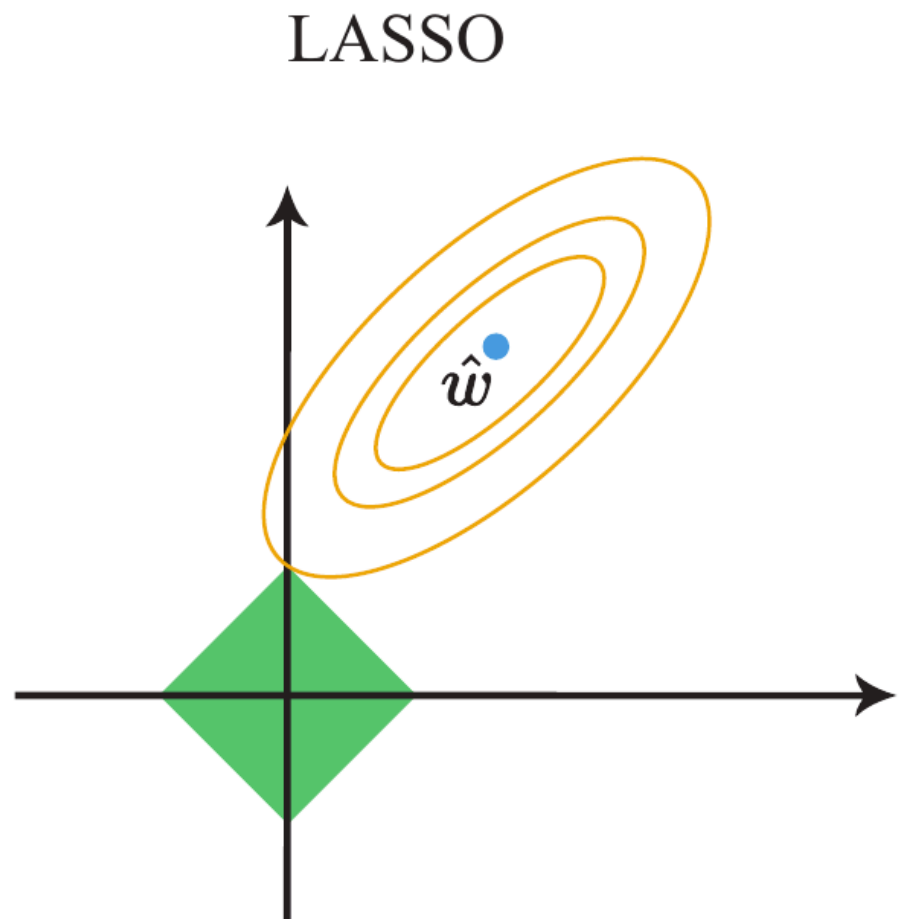
LASSO



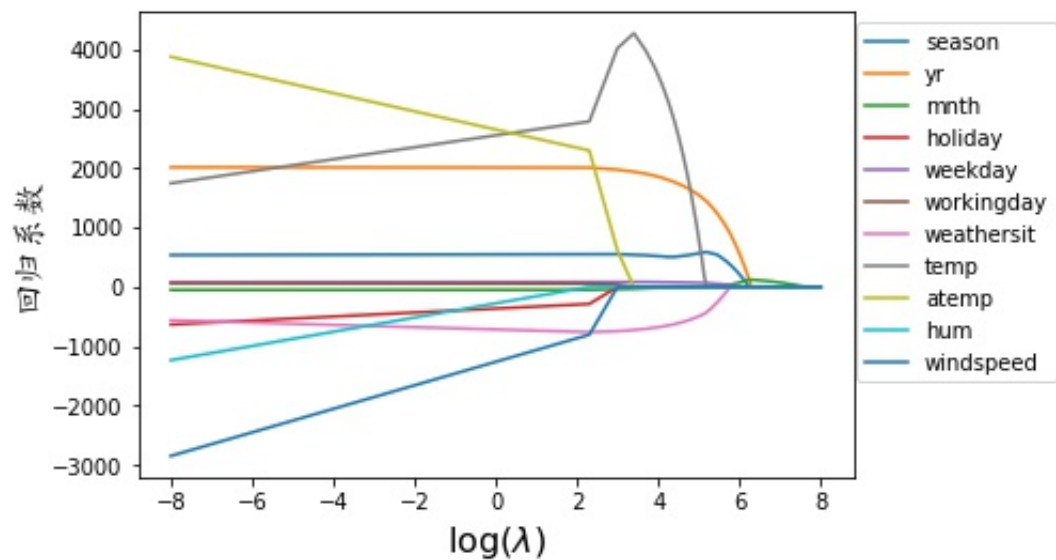
岭回归



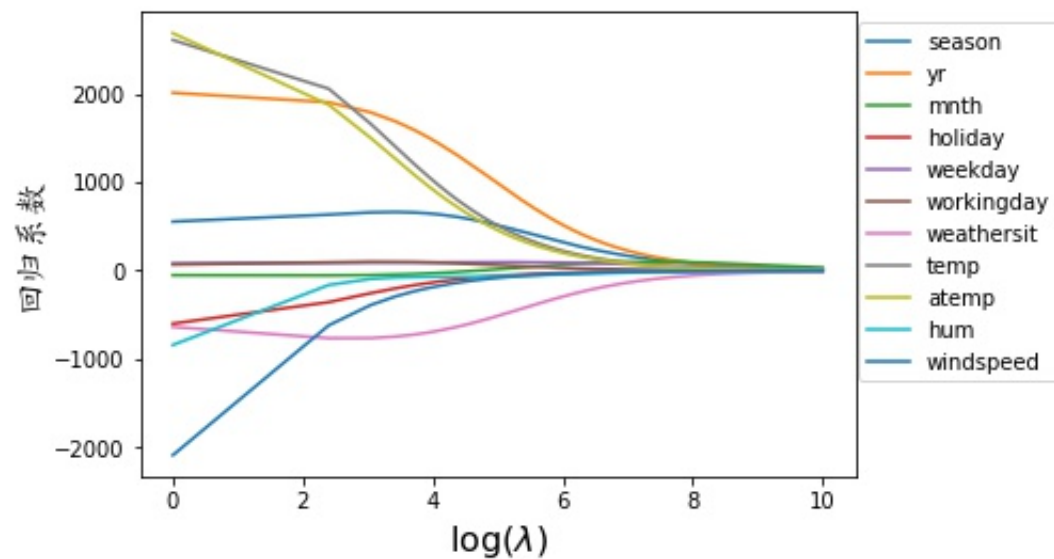
- 如右图所示：图中绿色区域表示约束区域，黄色线为残差平方和函数的等高线
- 通过添加 L1惩罚函数，LASSO 方法可以得到角点解，即稀疏的最优解  $\hat{\mathbf{w}}$ ，此时  $\hat{w}_i = 0$ ，我们可以将对应的变量从模型中删除。







LASSO

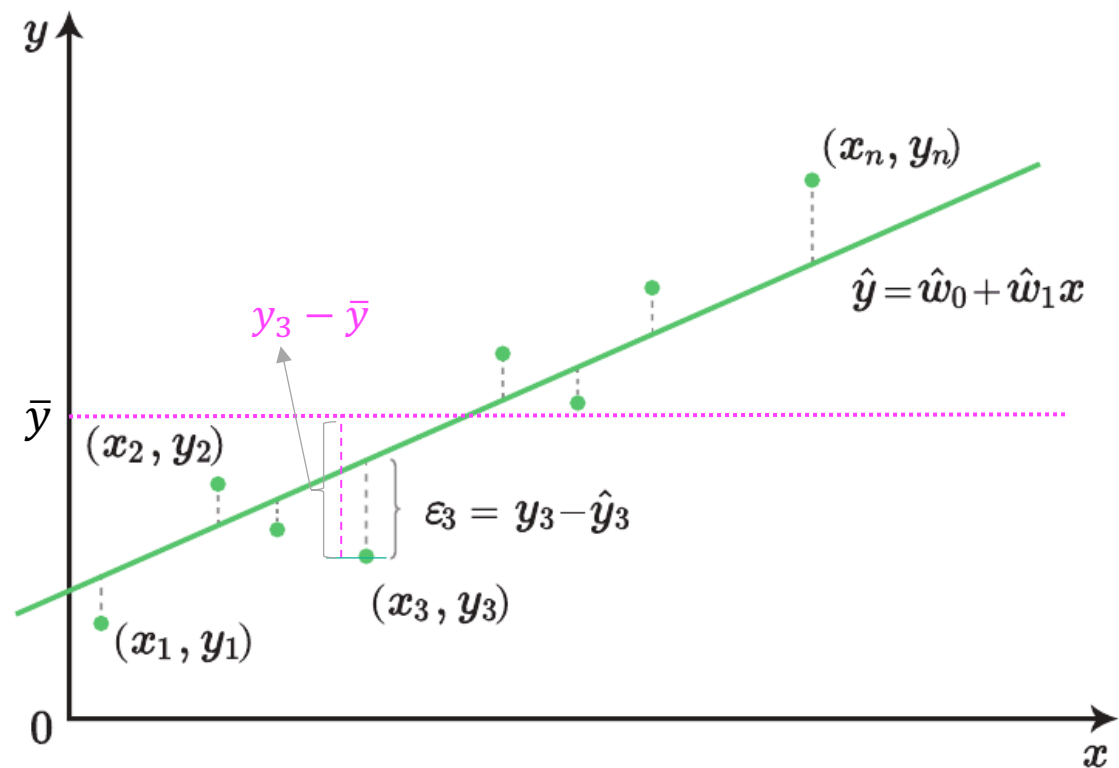


岭回归

随着  $\lambda$  增大，LASSO的特征系数逐个减小为0，可以做特征选择；而岭回归变量系数几乎同时趋近于0。

- 均方误差： $MSE(\mathbf{y}, \hat{\mathbf{y}}) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$
- 均方根误差： $RMSE(\mathbf{y}, \hat{\mathbf{y}}) = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$
- 平均绝对误差： $MAE(\mathbf{y}, \hat{\mathbf{y}}) = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$
- 决定系数： $R^2(\mathbf{y}, \hat{\mathbf{y}}) = 1 - \frac{SS_{res}}{SS_{tot}} = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$

其中 $y_i$ 为真实值， $\bar{y}$ 为真实值的平均值， $\hat{y}_i$ 为模型估计值



基于一份鲍鱼数据集，利用回归模型建立自动鲍鱼年龄预测模型。

1. 使用 Python 实现线性回归和岭回归算法，并与 Sklearn 中的实现进行对比
2. 借助 Sklearn 工具，对线性回归、岭回归和 LASSO 三种模型的预测效果使用 MAE 和决定系数进行效果评估
3. 残差图和正则化路径对模型表现进行分析





—— 数据酷客 ——



数据科学人工智能

- 回归：典型的有监督学习， $y$  为连续型

- 线性回归： $y = f(x) = w_1x_1 + w_2x_2 + \dots + w_dx_d + w_0$

- 目标函数： $L(\mathbf{w}) = \sum_{i=1}^n (\hat{y}_i - y_i)^2 = (\mathbf{X}\mathbf{w} - \mathbf{y})^T(\mathbf{X}\mathbf{w} - \mathbf{y})$

- 解： $\hat{\mathbf{w}} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$

- 岭回归：防止过度拟合

- 目标函数： $(\mathbf{X}\mathbf{w} - \mathbf{y})^T(\mathbf{X}\mathbf{w} - \mathbf{y}) + \lambda\|\mathbf{w}\|_2^2$

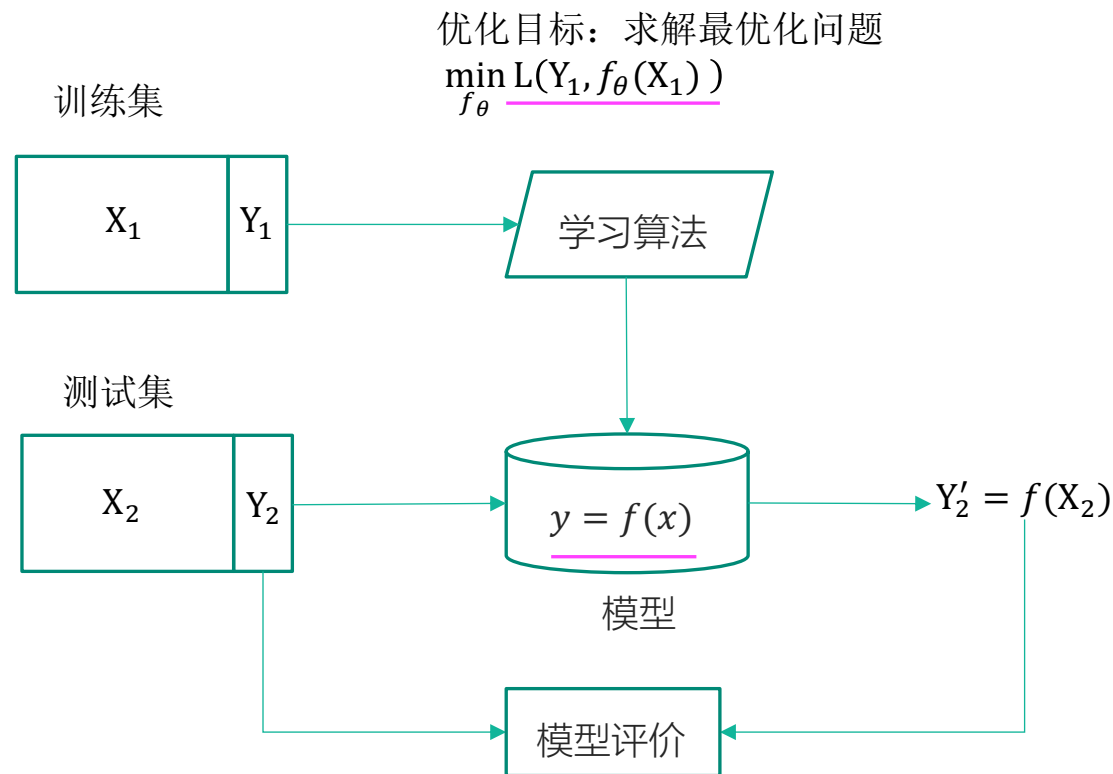
- 解： $\hat{\mathbf{w}}^{\text{ridge}} = (\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I})^{-1}\mathbf{X}^T\mathbf{y}$

- LASSO：同时进行特征选择，防止过度拟合

- 目标函数： $(\mathbf{X}\mathbf{w} - \mathbf{y})^T(\mathbf{X}\mathbf{w} - \mathbf{y}) + \lambda\|\mathbf{w}\|_1$

- 回归模型评价和分析

- MSE、RMSE、 $R^2$ 、残差图





—— 数据酷客 ——



数据科学人工智能



加入数据酷客交流群