

## 第二章 回归算法

在实际业务中经常预测一个值。“指标”，“价格”，“金额”，“分数”

↓  
高考成绩    房价    贷款、信用卡    考研分数

### 2.1 线性回归算法

提出者：英国生物学家、统计学家：“弗朗西斯·高尔顿”，父代与子代之间的身高关系。

线性回归：最简单、实用的算法之一。机器学习，掌握思路

假设有一场选秀大赛，组委会先了解某个选手的个人信息，例如：身高、体重、颜值，依据这些指标，决定一个选手是否能“晋级”，从而确定选手的“综合得分”组织如进行评估的？

是    否 → 淘汰    98分 → 晋级

选秀评估系统的建模过程：

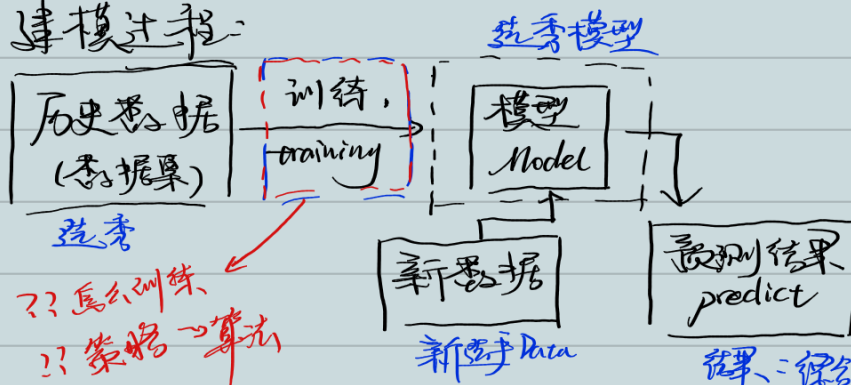
DataSet (数据集)

|                        | ID  | $X_1$<br>身高(m) | $X_2$<br>体重(kg) | $y$<br>综合得分 |
|------------------------|-----|----------------|-----------------|-------------|
| DataSet<br>↓<br>pandas | 001 | 1.73           | 60              | 82          |
|                        | 002 | 1.76           | 55              | 90          |
|                        | 003 | 1.82           | 65              | 87          |
|                        | 004 | 1.77           | 53              | 92          |
|                        | 005 | 1.71           | 65              | 75          |

$X_1 = \text{身高} \rightarrow$  “特征”、“属性”  
 $X_2 = \text{体重} \rightarrow$  “ ”  
}  $\Rightarrow$  综合得分 (标签)  
-----  $\rightarrow$  Label  
feature

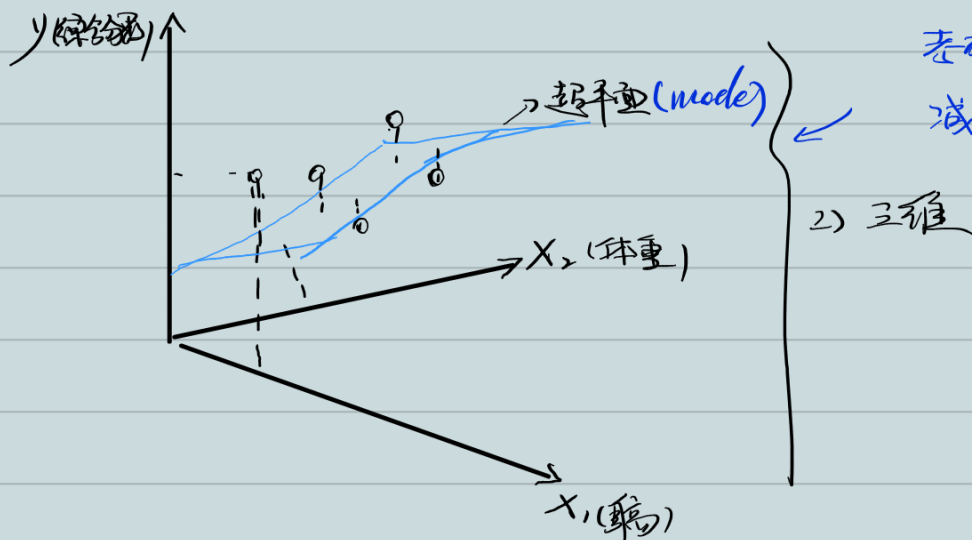
目的：得到一个公式描述  $x_1, x_2$  与  $y$  之间的关系。

机器学习建模过程：



## 2.2. 线性回归方程.

两个属性  $\rightarrow$  值 影响



圆点表示真实的得分值, 平面表示预测值, 可以观察到, 综合得分受  $X_1$  身高和  $X_2$  体重共同影响的. 每一个特征 ( $X_1, X_2$ ) 对  $y$  的影响有多大. 则对每个特征加入参数  $\rightarrow \theta$ . 公式:

$$\begin{cases} h_0(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 \\ \quad \quad \quad \text{偏差} \quad \quad \quad \text{身高} \quad \quad \quad \text{体重} \\ h_0(x) = \sum_{i=0}^n \theta_i x_i = \theta^T x \rightarrow \text{New} \rightarrow \star \rightarrow 95 \end{cases}$$

## 2.3. 误差项分析

