

# 统计方法与机器学习

## 第二章 线性回归分析

倪 蓓

Dase@ECNU

# 什么是线性回归分析？

- 在单因子方差分析模型中，所考虑的因子仅有  $a$  种不同的取值。

水平	观测到的响应变量			
1	$y_{11}$	$y_{12}$	...	$y_{1m}$
2	$y_{21}$	$y_{22}$	...	$y_{2m}$
$\vdots$	$\vdots$	$\vdots$		$\vdots$
$a$	$y_{a1}$	$y_{a2}$	...	$y_{am}$



因子	响应变量
1	$y_{11}$
1	$y_{12}$
$\vdots$	$\vdots$
1	$y_{1m}$
$\vdots$	$\vdots$
$a$	$y_{am}$

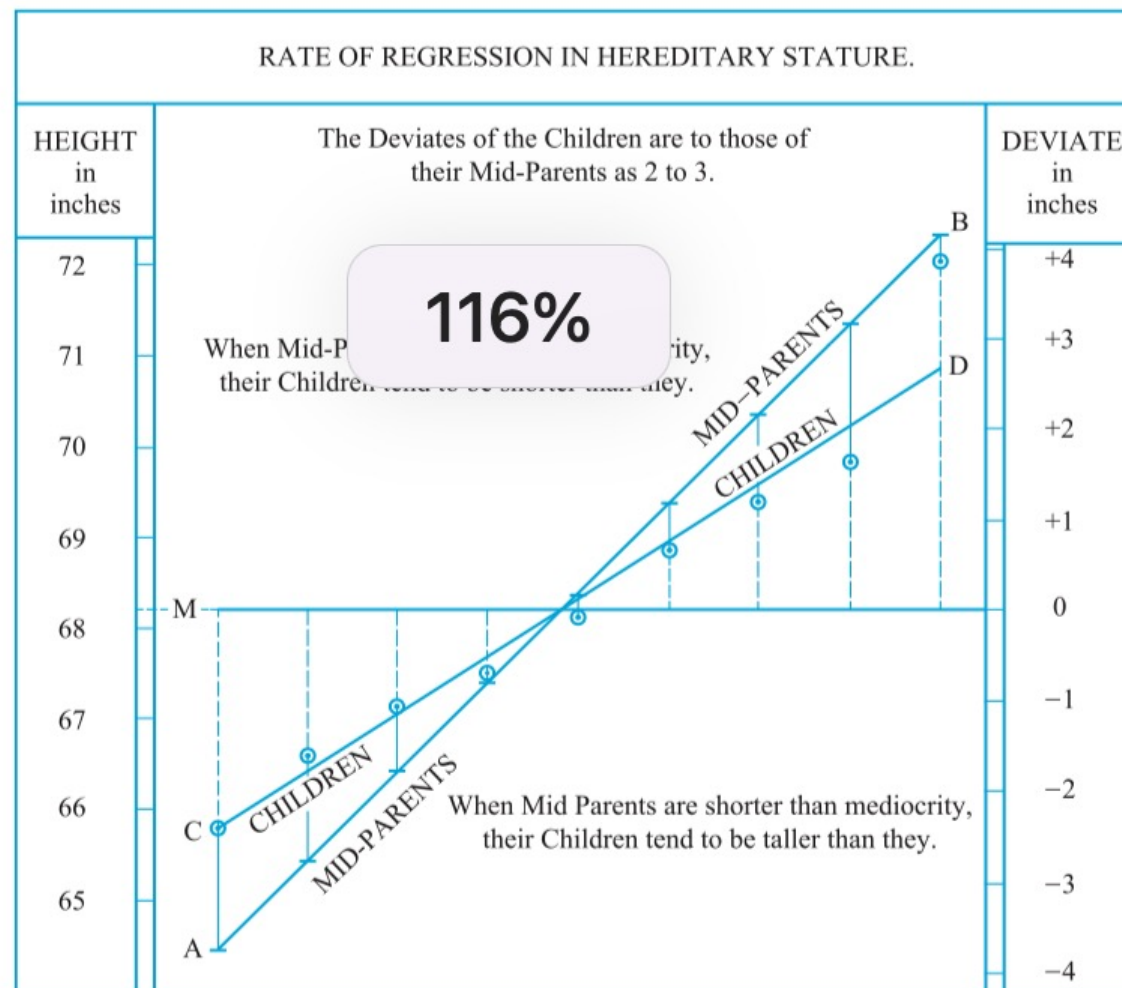
- 提问：每个样本中因子的取值都不相同，模型将会变成怎样？

自变量  
特征  
 $x$

响应变量  
标签  
 $y$

# 什么是线性回归分析？

- 回归源于生物遗传问题的研究。



# 什么是线性回归分析？

- **回归**的本质是构建响应变量  $y$  与自变量  $x$  的一种统计关系，即

$$f(x) = E(y|x)$$

- 特例：

- $f(x) = \beta_0 + \beta_1 x$
- $f(x) = \max(0, \beta_0 + \beta_1 x)$
- $f(x) = g_1 \circ g_2 \circ \cdots \circ g_n(\beta_0 + \beta_1 x)$

- 应用场景：

- 信用卡授信金额与申请人年龄、职业、收入、消费偏好有关。
- 房屋价格与房屋面积、卧室数量、地段、建造时间、交通便利有关。
- 景区的游览人数与天气情况、是否为节假日、景区类型等有关。

# 目录

## 1. 多元线性回归 ( Multiple Linear Regression )

I. 模型

II. 参数估计

III. 中心化与标准化

IV. 显著性检验

V. 估计与预测

# 目录

## 1. 多元线性回归 ( Multiple Linear Regression )

I. 模型

II. 参数估计

III. 中心化与标准化

IV. 显著性检验

V. 估计与预测

# 定义

- 线性回归模型为

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_p x_p + \varepsilon$$

- 称  $y$  为响应变量或标签，通常认为是一个连续型随机变量。
- 称  $x_j$  为第  $j$  个自变量或特征，通常认为是确定性的变量。
- 称  $\beta_j$  为回归系数，共有  $p + 1$  个未知参数。
- 称  $\varepsilon$  为随机误差，并假定

$$E(\varepsilon) = 0,$$

$$\text{Var}(\varepsilon) = \sigma^2.$$

# 数据

- 共有  $n$  组观测数据

$$\{(x_{i1}, x_{i2}, \dots, x_{ip}, y_i) : i = 1, 2, \dots, n\}$$

- 于是有

$$\begin{cases} y_1 = \beta_0 + \beta_1 x_{11} + \beta_2 x_{12} + \dots + \beta_p x_{1p} + \varepsilon_1 \\ \vdots \\ y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + \varepsilon_i \\ \vdots \\ y_n = \beta_0 + \beta_1 x_{n1} + \beta_2 x_{n2} + \dots + \beta_p x_{np} + \varepsilon_n \end{cases}$$



## 数据

$$\mathbf{y} = (y_1, y_2, \dots, y_n)^\top$$

$$\boldsymbol{\beta} = (\beta_0, \beta_1, \beta_2, \dots, \beta_p)^\top$$

$$\boldsymbol{\varepsilon} = (\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n)^\top$$

$$\begin{cases} y_1 = \beta_0 + \beta_1 x_{11} + \beta_2 x_{12} + \dots + \beta_p x_{1p} + \varepsilon_1 \\ y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + \varepsilon_i \\ y_n = \beta_0 + \beta_1 x_{n1} + \beta_2 x_{n2} + \dots + \beta_p x_{np} + \varepsilon_n \end{cases}$$

- 矩阵形式为

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

$$\mathbf{X} = \begin{pmatrix} 1 & x_{11} & x_{12} & \dots & x_{1p} \\ 1 & x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{np} \end{pmatrix}$$

# 基本假定

- 设计矩阵  $\mathbf{X}$  是确定性的，且  $\text{rank}(\mathbf{X}) = p + 1 < n$ 。
  - 在本课程中， $\mathbf{X}$  是确定性的，而非随机性的；
  - $\mathbf{X}$  是一个列满秩矩阵。
- 随机误差  $\varepsilon$  是零均值且等方差的，即
  - $E(\varepsilon_i) = 0, i = 1, 2, \dots, n$
  - $\text{Cov}(\varepsilon_i, \varepsilon_j) = \begin{cases} \sigma^2, & i = j \\ 0, & i \neq j \end{cases}$
  - 称为**高斯-马尔可夫条件**。

# 基本假定

- 进一步，假定随机误差服从正态分布，即

$$\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

- 期望： $E(\boldsymbol{\varepsilon}) = \mathbf{0}$ ；
- 方差-协方差矩阵： $\text{Var}(\boldsymbol{\varepsilon}) = \sigma^2 \mathbf{I}_n$ 。
- 因为

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

所以，

$$\mathbf{y} \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n)$$

- 期望： $E(\mathbf{y}) = \mathbf{X}\boldsymbol{\beta}$ ；
- 方差-协方差矩阵： $\text{Var}(\mathbf{y}) = \sigma^2 \mathbf{I}_n$ 。

# 总结

- 多元线性回归模型为

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

- 基本假定1： $\mathbf{X}$  是确定性的，且  $\text{rank}(\mathbf{X})=p+1 < n$ 。
- 基本假定2： $\boldsymbol{\varepsilon}$  是零均值且等方差。
- 基本假定3： $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_n)$

# 目录

## 1. 多元线性回归 ( Multiple Linear Regression )

I. 模型

II. 参数估计

III. 中心化与标准化

IV. 显著性检验

V. 估计与预测

# 问题定义

- 多元线性回归模型为

$$y = X\beta + \varepsilon$$

- 待估计参数

- 回归系数  $\beta = (\beta_0, \beta_1, \beta_2 \cdots, \beta_p)^\top$  ;

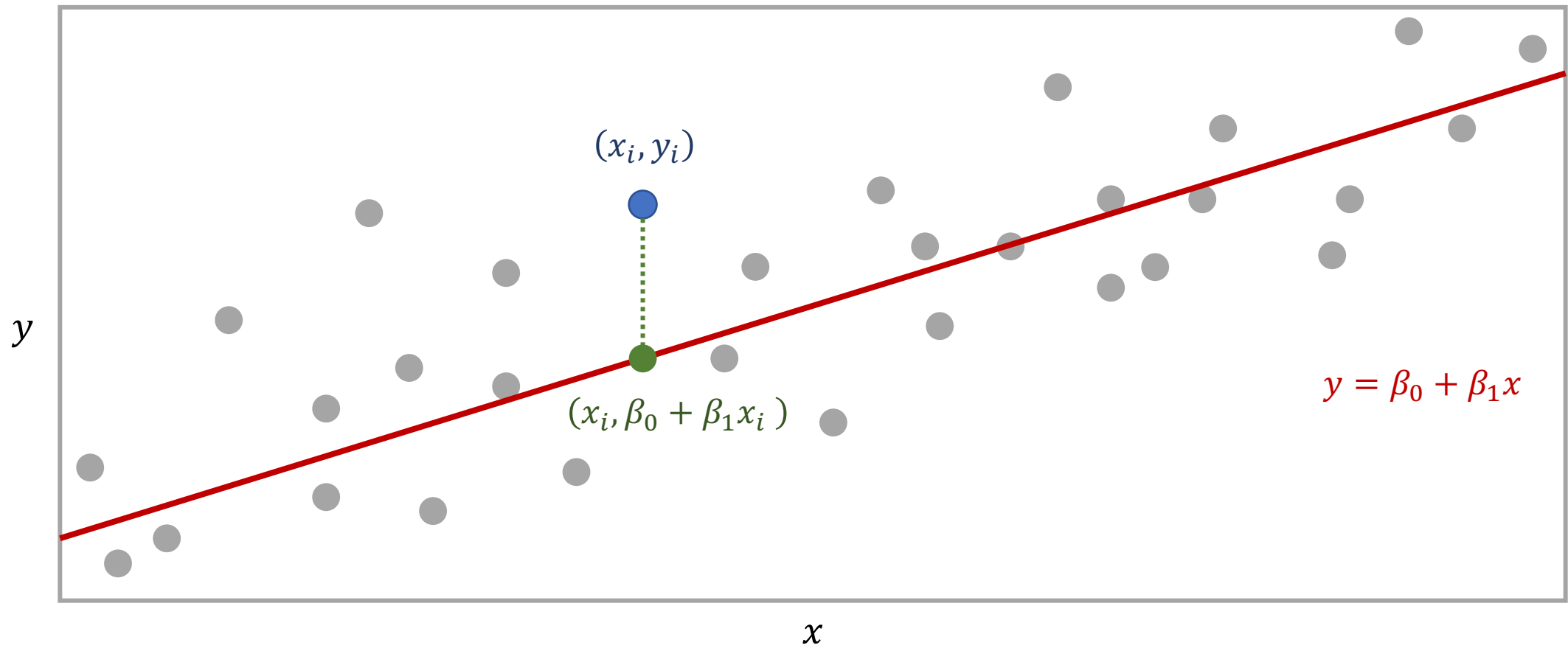
- 误差方差  $\sigma^2$  ;

- 估计方法

- 最小二乘估计
- 极大似然估计

# 最小二乘估计

• 基本思想：拟合



# 最小二乘估计

- 基本思想：拟合
- 线性回归模型

$$E(y|\mathbf{x}) = \mathbf{x}^\top \boldsymbol{\beta} = \beta_0 + \beta_1 x_1 + \cdots + \beta_p x_p$$

- 第  $i$  个数据点的实际观测值为  $(\mathbf{x}_i, y_i)$
- 第  $i$  个数据点的拟合值为  $(\mathbf{x}_i, \mathbf{x}_i^\top \boldsymbol{\beta})$
- 差异为

$$y_i - \mathbf{x}_i^\top \boldsymbol{\beta}$$



# 最小二乘估计

- 损失函数定义为

$$Q(\boldsymbol{\beta}) = \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2$$

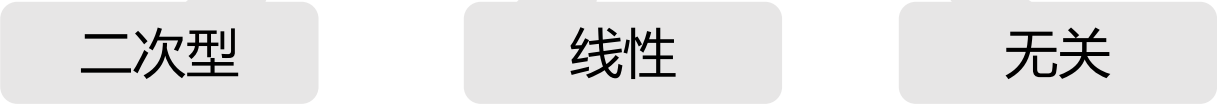
- 最小二乘估计为

$$\hat{\boldsymbol{\beta}}_{\text{LS}} = \operatorname{argmin}_{\boldsymbol{\beta}} Q(\boldsymbol{\beta})$$

# 最小二乘估计

- 求解过程
  - 损失函数的另一种形式

$$\begin{aligned} Q(\boldsymbol{\beta}) &= \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2 \\ &= (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \\ &= (\mathbf{X}\boldsymbol{\beta})^\top (\mathbf{X}\boldsymbol{\beta}) - 2(\mathbf{X}\boldsymbol{\beta})^\top \mathbf{y} + \mathbf{y}^\top \mathbf{y} \\ &= \boxed{\boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X} \boldsymbol{\beta}} - \boxed{2\boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{y}} + \boxed{\mathbf{y}^\top \mathbf{y}} \end{aligned}$$



# 最小二乘估计

- 求解过程
  - 求导，即

$$\frac{\partial Q(\beta)}{\partial \beta} = \frac{\partial (\beta^T X^T X \beta - 2\beta^T X^T y + y^T y)}{\partial \beta} = 2X^T X \beta - 2X^T y$$

- 令  $\frac{\partial Q(\beta)}{\partial \beta} = 0$ ，可得

$$X^T X \beta = X^T y$$

- 根据基本假设1可知， $X^T X$  是满秩的。
- 最小二乘估计为

$$\hat{\beta}_{LS} = (X^T X)^{-1} X^T y$$

## 知识回顾——向量求导

假定

- $x$  是一个  $p$  维向量
- $a$  是一个  $p$  维向量
- $B$  是  $p \times p$  矩阵

结论

- 线性：

$$\frac{\partial (x^T a)}{\partial x} = a$$

- 二次型：

$$\frac{\partial (x^T B x)}{\partial x} = (B + B^T)x$$

# 最小二乘估计

- 最小二乘估计为

$$\hat{\beta}_{LS} = (X^T X)^{-1} X^T y$$

- 反思：为什么 基本假定1： $\text{rank}(X) = p + 1 < n$  是重要的？
- 在最小二乘估计中，矩阵  $X^T X$  需要求逆矩阵。
- 因为  $\text{rank}(X^T X) \leq \text{rank}(X)$ ，基本假定1能够保证  $\text{rank}(X^T X) = p + 1$ 。
- 基本假定1的直观解释
  - 特征之间不（完全）线性相关；
  - 样本量需要（远）大于特征个数。

# 最小二乘估计

- 最小二乘估计为

$$\hat{\beta}_{LS} = (X^T X)^{-1} X^T y$$

- 反思：为什么 基本假定1： $\text{rank}(X) = p + 1 < n$  是重要的？
- 在最小二乘估计中，矩阵  $X^T X$  需要求逆矩阵。
- 因为  $\text{rank}(X^T X) \leq \text{rank}(X)$ ，基本假定1能够保证  $\text{rank}(X^T X) = p + 1$ 。
- 基本假定1的直观解释
  - 特征之间不（完全）线性相关；
  - 样本量需要（远）大于特征个数。

# 极大似然估计

- 基本思想：似然
- 极大似然估计依赖于数据分布假定，即

$$\mathbf{y} \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n)$$

- $\mathbf{y}$  的联合密度函数为

$$p(\mathbf{y}; \boldsymbol{\beta}, \sigma^2) = (2\pi)^{-n/2} |\sigma^2 \mathbf{I}_n|^{-1/2} \exp \left\{ -\frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\sigma^2 \mathbf{I}_n)^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}$$

- $(\boldsymbol{\beta}, \sigma^2)$  的似然函数为

$$L(\boldsymbol{\beta}, \sigma^2) = (2\pi)^{-p/2} |\sigma^2 \mathbf{I}_n|^{-1/2} \exp \left\{ -\frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\sigma^2 \mathbf{I}_n)^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}$$

# 极大似然估计

- 基本思想：似然
- $(\boldsymbol{\beta}, \sigma^2)$  的似然函数为

$$L(\boldsymbol{\beta}, \sigma^2) = (2\pi)^{-n/2} |\sigma^2 \mathbf{I}_n|^{-1/2} \exp \left\{ -\frac{1}{2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\sigma^2 \mathbf{I}_n)^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}$$

$$|\sigma^2 \mathbf{I}_n| = \begin{vmatrix} \sigma^2 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \sigma^2 \end{vmatrix} = (\sigma^2)^n$$

$$(\sigma^2 \mathbf{I}_n)^{-1} = \begin{pmatrix} \sigma^2 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \sigma^2 \end{pmatrix}^{-1} = \sigma^{-2} \begin{pmatrix} 1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & 1 \end{pmatrix}$$

- 极大似然估计为

$$\begin{aligned} (\hat{\boldsymbol{\beta}}_{\text{MLE}}, \hat{\sigma}_{\text{MLE}}^2) &= \operatorname{argmax} L(\boldsymbol{\beta}, \sigma^2) \\ &= \operatorname{argmax} \ln(L(\boldsymbol{\beta}, \sigma^2)) \end{aligned}$$

# 极大似然估计

- 求解过程
  - 对数似然函数为

$$\ln(L(\boldsymbol{\beta}, \sigma^2)) = -\frac{n}{2}\ln(2\pi) - \frac{n}{2}\ln \sigma^2 - \frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

- 对  $\boldsymbol{\beta}$  和  $\sigma^2$  求偏导，即

$$\begin{cases} \frac{\partial \ln(L(\boldsymbol{\beta}, \sigma^2))}{\partial \boldsymbol{\beta}} = -\frac{1}{\sigma^2}(\mathbf{X}^\top \mathbf{X}\boldsymbol{\beta} - \mathbf{X}^\top \mathbf{y}) \\ \frac{\partial \ln(L(\boldsymbol{\beta}, \sigma^2))}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{n}{2(\sigma^2)^2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \end{cases}$$



# 极大似然估计

- 求解过程

- 极大似然估计为

$$\begin{cases} \hat{\boldsymbol{\beta}}_{\text{MLE}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \\ \hat{\sigma}^2_{\text{MLE}} = \frac{1}{n} (\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}}_{\text{MLE}})^\top (\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}}_{\text{MLE}}) \end{cases}$$

- 说明

- $\boldsymbol{\beta}$  的最小二乘估计与极大似然估计是一致的，记  $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$ 。
- $\hat{\sigma}^2_{\text{MLE}}$  不是无偏估计，但是相合估计。

# 参数估计的性质

## 定理2.5

在多元线性回归模型中， $\beta$  的估计为

$$\hat{\beta} = (X^T X)^{-1} X^T y$$

有以下两个重要性质：

- $\hat{\beta}$  的期望为

$$E(\hat{\beta}) = \beta$$

- $\hat{\beta}$  的方差-协方差矩阵为

$$\text{Var}(\hat{\beta}) = \sigma^2 (X^T X)^{-1}$$

# 参数估计的性质

**定理2.5** 在多元线性回归模型中， $\beta$  的估计为

$$\hat{\beta} = (X^T X)^{-1} X^T y$$

有以下两个重要性质：

- $\hat{\beta}$  的期望为

$$E(\hat{\beta}) = \beta$$

- $\hat{\beta}$  的方差-协方差矩阵为

$$\text{Var}(\hat{\beta}) = \sigma^2 (X^T X)^{-1}$$

# 参数估计的性质

## 定理2.5的证明

$$\begin{aligned} E(\hat{\beta}) &= E((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}) \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top E(\mathbf{y}) \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top E(\mathbf{X}\beta + \varepsilon) \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{X}\beta + E(\varepsilon)) \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}\beta = \beta \end{aligned}$$

知识回顾

因为

$$\mathbf{y} = \mathbf{X}\beta + \varepsilon$$

所以，

- 期望： $E(\mathbf{y}) = \mathbf{X}\beta$ ；
- 方差-协方差矩阵： $\text{Var}(\mathbf{y}) = \sigma^2 \mathbf{I}_n$ 。

# 参数估计的性质

## 定理2.5

在多元线性回归模型中， $\beta$  的估计为

$$\hat{\beta} = (X^T X)^{-1} X^T y$$

有以下两个重要性质：

- $\hat{\beta}$  的期望为

$$E(\hat{\beta}) = \beta$$

- $\hat{\beta}$  的方差-协方差矩阵为

$$\text{Var}(\hat{\beta}) = \sigma^2 (X^T X)^{-1}$$

# 参数估计的性质

## 定理2.5的证明

$$\begin{aligned}\text{Var}(\hat{\beta}) &= \text{Var}((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}) \\&= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \text{Var}(\mathbf{y}) ((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top)^\top \\&= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \text{Var}(\mathbf{X}\beta + \varepsilon) ((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top)^\top \\&= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\sigma^2 \mathbf{I}_n) \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \\&= \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \\&= \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}\end{aligned}$$

知识回顾

因为

$$\mathbf{y} = \mathbf{X}\beta + \varepsilon$$

所以，

- 期望： $E(\mathbf{y}) = \mathbf{X}\beta$ ；
- 方差-协方差矩阵： $\text{Var}(\mathbf{y}) = \sigma^2 \mathbf{I}_n$ 。

# 参数估计的几何解释

- $y$  的拟合值为

$$\begin{aligned}\hat{y} &= X\hat{\beta} \\ &= X(X^T X)^{-1} X^T y\end{aligned}$$

- 定义

$$H = X(X^T X)^{-1} X^T$$

为**帽子矩阵** ( hat matrix ) 。

- 于是 ,

$$\hat{y} = Hy$$

# 参数估计的几何解释

## 命题2.1 帽子矩阵

$$\boldsymbol{H} = \boldsymbol{X}(\boldsymbol{X}^\top \boldsymbol{X})^{-1} \boldsymbol{X}^\top$$

有以下三个重要性质：

- $\boldsymbol{H}$  是  $n$  阶对称矩阵，即  $\boldsymbol{H} = \boldsymbol{H}^\top$ ；
- $\boldsymbol{H}$  是幂等矩阵，即  $\boldsymbol{H} = \boldsymbol{H}^2$ ；
- $\boldsymbol{H}$  的迹为  $p + 1$ ，即  $\text{tr}(\boldsymbol{H}) = p + 1$ 。



# 参数估计的几何解释

## 命题2.1

帽子矩阵

$$H = X(X^T X)^{-1} X^T$$

有以下三个重要性质：

- $H$  是  $n$  阶对称矩阵，即  $H = H^T$ ；
- $H$  是幂等矩阵，即  $H = H^2$ ；
- $H$  的迹为  $p + 1$ ，即  $\text{tr}(H) = p + 1$ 。

# 参数估计的性质

## 命题2.1的证明

$$\begin{aligned} \boldsymbol{H}^\top &= (\boldsymbol{X}(\boldsymbol{X}^\top \boldsymbol{X})^{-1} \boldsymbol{X}^\top)^\top \\ &= (\boldsymbol{X}^\top)^\top ((\boldsymbol{X}^\top \boldsymbol{X})^{-1})^\top \boldsymbol{X}^\top \\ &= \boldsymbol{X}(\boldsymbol{X}^\top \boldsymbol{X})^{-1} \boldsymbol{X}^\top = \boldsymbol{H} \end{aligned}$$

# 参数估计的几何解释

## 命题2.1

帽子矩阵

$$H = X(X^T X)^{-1} X^T$$

有以下三个重要性质：

- $H$  是  $n$  阶对称矩阵，即  $H = H^T$ ；
- $H$  是幂等矩阵，即  $H = H^2$ ；
- $H$  的迹为  $p + 1$ ，即  $\text{tr}(H) = p + 1$ 。

# 参数估计的性质

## 命题2.1的证明

$$\begin{aligned} H^2 &= (\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top)^2 \\ &= (\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) (\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) \\ &= \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} (\mathbf{X}^\top \mathbf{X}) (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \\ &= \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top = H \end{aligned}$$

# 参数估计的几何解释

## 命题2.1

帽子矩阵

$$\boldsymbol{H} = \boldsymbol{X}(\boldsymbol{X}^\top \boldsymbol{X})^{-1} \boldsymbol{X}^\top$$

有以下三个重要性质：

- $\boldsymbol{H}$  是  $n$  阶对称矩阵，即  $\boldsymbol{H} = \boldsymbol{H}^\top$ ；
- $\boldsymbol{H}$  是幂等矩阵，即  $\boldsymbol{H} = \boldsymbol{H}^2$ ；
- $\boldsymbol{H}$  的迹为  $p + 1$ ，即  $\text{tr}(\boldsymbol{H}) = p + 1$ 。

# 参数估计的性质

## 命题2.1的证明

$$\begin{aligned}\text{tr}(\mathbf{H}) &= \text{tr}(\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) \\ &= \text{tr}((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}) \\ &= \text{tr}(\mathbf{I}_{p+1}) = p + 1\end{aligned}$$

# 参数估计的几何解释

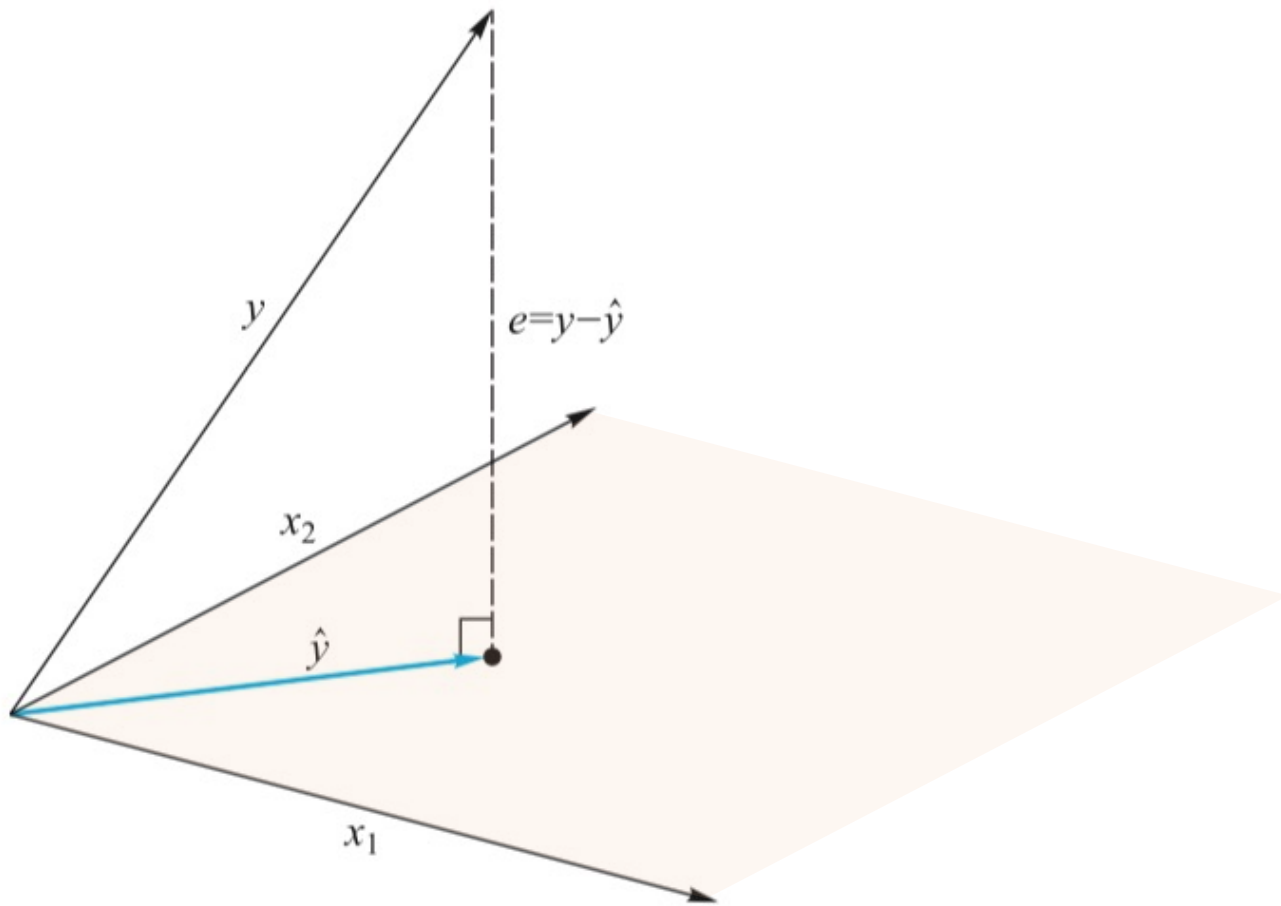
- 残差定义为

$$\begin{aligned} \mathbf{e} &= \mathbf{y} - \hat{\mathbf{y}} \\ &= \mathbf{y} - \mathbf{H}\mathbf{y} \\ &= (\mathbf{I} - \mathbf{H})\mathbf{y} \end{aligned}$$

- 几何关系：拟合值  $\hat{\mathbf{y}}$  与残差  $\mathbf{e}$  垂直，即

$$\begin{aligned} \hat{\mathbf{y}}^\top \mathbf{e} &= (\mathbf{H}\mathbf{y})^\top (\mathbf{I} - \mathbf{H})\mathbf{y} \\ &= \mathbf{y}^\top \mathbf{H}^\top (\mathbf{I} - \mathbf{H})\mathbf{y} = 0 \end{aligned}$$

# 参数估计的几何解释





# 参数估计的几何解释

- 残差的性质
- 残差的期望

$$E(\mathbf{e}) = E((\mathbf{I} - \mathbf{H})\mathbf{y}) = (\mathbf{I} - \mathbf{H})E(\mathbf{y}) = (\mathbf{I} - \mathbf{H})\mathbf{X}\boldsymbol{\beta} = \mathbf{0}$$

- 残差的方差-协方差矩阵

$$\begin{aligned}\text{Var}(\mathbf{e}) &= \text{Cov}(\mathbf{e}, \mathbf{e}) \\ &= \text{Cov}((\mathbf{I} - \mathbf{H})\mathbf{y}, (\mathbf{I} - \mathbf{H})\mathbf{y}) \\ &= (\mathbf{I} - \mathbf{H})\text{Cov}(\mathbf{y}, \mathbf{y})(\mathbf{I} - \mathbf{H})^\top \\ &= \sigma^2(\mathbf{I} - \mathbf{H})\mathbf{I}_n(\mathbf{I} - \mathbf{H})^\top \\ &= \sigma^2(\mathbf{I} - \mathbf{H})\end{aligned}$$

# 参数估计的几何解释

- 提问：能否得到  $\sigma^2$  的无偏估计？
- 一般而言，采用

$$\hat{\sigma}^2 = \frac{1}{n - (p + 1)} \mathbf{e}^\top \mathbf{e}$$

作为  $\sigma^2$  的无偏估计。

# 总结

- 回归系数  $\beta$  的估计为

$$\hat{\beta} = (X^T X)^{-1} X^T y$$

- 随机误差的方差  $\sigma^2$  的估计
  - 无偏估计

$$\hat{\sigma}^2 = \frac{1}{n - (p + 1)} e^T e$$

- 极大似然估计

$$\hat{\sigma}_{\text{MLE}}^2 = \frac{1}{n} e^T e$$

# 目录

## 1. 多元线性回归 ( Multiple Linear Regression )

I. 模型

II. 参数估计

III. 中心化与标准化

IV. 显著性检验

V. 估计与预测

# 回顾：最小二乘估计

- 多元线性回归模型

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_p x_p + \varepsilon$$

- 现有数据  $\{(x_{i1}, x_{i2}, \cdots, x_{ip}, y_i): i = 1, 2, \cdots, n\}$
- 可以得到回归系数的估计

$$\begin{aligned}\hat{\boldsymbol{\beta}} &= (\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2, \cdots, \hat{\beta}_p)^\top \\ &= (\hat{\beta}_0, \hat{\boldsymbol{\beta}}_1^\top)^\top\end{aligned}$$

- 定义

$$\mathbf{y} = (y_1, \cdots, y_n)^\top, \mathbf{X} = (\mathbf{1}_n, \mathbf{x}_1, \mathbf{x}_2, \cdots, \mathbf{x}_p) = (\mathbf{1}_n, \mathbf{X}_0)$$

# 回顾：最小二乘估计

• 回归系数的估计为

$$\begin{aligned}
 \hat{\beta} &= (X^T X)^{-1} X^T y \\
 &= ((\mathbf{1}_n, \mathbf{X}_o)^T (\mathbf{1}_n, \mathbf{X}_o))^{-1} (\mathbf{1}_n, \mathbf{X}_o)^T y \\
 &= \begin{pmatrix} \mathbf{1}_n^T \mathbf{1}_n & \mathbf{1}_n^T \mathbf{X}_o \\ \mathbf{X}_o^T \mathbf{1}_n & \mathbf{X}_o^T \mathbf{X}_o \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{1}_n^T \\ \mathbf{X}_o^T \end{pmatrix} y \\
 &= \begin{pmatrix} n^{-1} + n^{-2} \mathbf{1}_n^T \mathbf{X}_o \mathbf{A}_o \mathbf{X}_o^T \mathbf{1}_n & -n^{-1} \mathbf{1}_n^T \mathbf{X}_o \mathbf{A}_o \\ -n^{-1} \mathbf{A}_o \mathbf{X}_o^T \mathbf{1}_n & \mathbf{A}_o \end{pmatrix} \begin{pmatrix} \mathbf{1}_n^T \\ \mathbf{X}_o^T \end{pmatrix} y \\
 &= \begin{pmatrix} n^{-1} \mathbf{1}_n^T + n^{-2} \mathbf{1}_n^T \mathbf{X}_o \mathbf{A}_o \mathbf{X}_o^T \mathbf{1}_n \mathbf{1}_n^T - n^{-1} \mathbf{1}_n^T \mathbf{X}_o \mathbf{A}_o \mathbf{X}_o^T \\ -n^{-1} \mathbf{A}_o \mathbf{X}_o^T \mathbf{1}_n \mathbf{1}_n^T + \mathbf{A}_o \mathbf{X}_o^T \end{pmatrix} y
 \end{aligned}$$

$$\begin{aligned}
 &\mathbf{1}_n^T \mathbf{1}_n \\
 &= (1 \quad \dots \quad 1) \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \\
 &= n
 \end{aligned}$$

其中， $\mathbf{A}_o = (\mathbf{X}_o^T \mathbf{X}_o - n^{-1} \mathbf{X}_o^T \mathbf{1}_n \mathbf{1}_n^T \mathbf{X}_o)^{-1}$

知识回顾——分块矩阵求逆

$A$  是可逆矩阵，如果  $D - CA^{-1}B$  可逆，那么有

$$\begin{pmatrix} A & B \\ C & D \end{pmatrix}^{-1} = \begin{pmatrix} E_{11} & E_{12} \\ E_{21} & E_{22} \end{pmatrix}$$

- $E_{11} = A^{-1} + A^{-1} B E_{22} C A^{-1}$
- $E_{12} = -A^{-1} B E_{22}$
- $E_{21} = -E_{22} C A^{-1}$
- $E_{22} = (D - C A^{-1} B)^{-1}$

问题定义

		原始数据	中心化	标准化
数据变换	响应变量	$y = (y_1, \cdots, y_n)^\top$	$\mathbf{y}^* = (y_1^*, \cdots, y_n^*)^\top$ $y_i^* = y_i - \bar{y}, \bar{y} = n^{-1} \sum_{i=1}^n y_i$	$\mathbf{y}^{**} = (y_1^{**}, \cdots, y_n^{**})^\top$ $y_i^{**} = \frac{y_i^*}{\sqrt{l_{yy}}}, l_{yy} = \sum_{i=1}^n (y_i^*)^2$
	自变量	$\mathbf{X} = (\mathbf{1}_n, \mathbf{X}_o),$ $\mathbf{X}_o = \{x_{ij}\}_{n \times p}$	$\mathbf{X}^* = (\mathbf{1}_n, \mathbf{X}_c), \mathbf{X}_c = \{x_{ij}^*\}_{n \times p}$ $x_{ij}^* = x_{ij} - \bar{x}_j, \bar{x}_j = n^{-1} \sum_{i=1}^n x_{ij}$	$\mathbf{X}^{**} = (\mathbf{1}_n, \mathbf{X}_s), \mathbf{X}_s = \{x_{ij}^{**}\}_{n \times p}$ $x_{ij}^{**} = \frac{x_{ij}^*}{\sqrt{l_{jj}}}, l_{jj} = \sum_{i=1}^n (x_{ij}^*)^2$
参数估计	$\beta_0$	$\hat{\beta}_0$	$\hat{\beta}_{c,0} = ?$	$\hat{\beta}_{s,0} = ?$
	$\begin{pmatrix} \beta_1 \\ \vdots \\ \beta_p \end{pmatrix}$	$\hat{\beta}_1$	$\hat{\beta}_{c,1} = ?$	$\hat{\beta}_{s,1} = ?$

# 中心化

• 回归系数的估计为

$$\begin{aligned}
 \hat{\beta}_c &= ((X^*)^\top X^*)^{-1} (X^*)^\top y^* \\
 &= ((\mathbf{1}_n, X_c)^\top (\mathbf{1}_n, X_c))^{-1} (\mathbf{1}_n, X_c)^\top y^* \\
 &= \begin{pmatrix} \mathbf{1}_n^\top \mathbf{1}_n & \mathbf{1}_n^\top X_c \\ X_c^\top \mathbf{1}_n & X_c^\top X_c \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{1}_n^\top \\ X_c^\top \end{pmatrix} y^* \\
 &= \begin{pmatrix} n^{-1} + n^{-2} \mathbf{1}_n^\top X_c A_c X_c^\top \mathbf{1}_n & -n^{-1} \mathbf{1}_n^\top X_c A_c \\ -n^{-1} A_c X_c^\top \mathbf{1}_n & A_c \end{pmatrix} \begin{pmatrix} \mathbf{1}_n^\top \\ X_c^\top \end{pmatrix} y^* \\
 &= \begin{pmatrix} n^{-1} \mathbf{1}_n^\top + n^{-2} \mathbf{1}_n^\top X_c A_c X_c^\top \mathbf{1}_n \mathbf{1}_n^\top - n^{-1} \mathbf{1}_n^\top X_c A_c X_c^\top \\ -n^{-1} A_c X_c^\top \mathbf{1}_n \mathbf{1}_n^\top + A_c X_c^\top \end{pmatrix} y^*
 \end{aligned}$$

其中,  $A_c = (X_c^\top X_c - n^{-1} X_c^\top \mathbf{1}_n \mathbf{1}_n^\top X_c)^{-1}$

## 知识回顾——分块矩阵求逆

$A$  是可逆矩阵, 如果  $D - CA^{-1}B$  可逆, 那么有

$$\begin{pmatrix} A & B \\ C & D \end{pmatrix}^{-1} = \begin{pmatrix} E_{11} & E_{12} \\ E_{21} & E_{22} \end{pmatrix}$$

- $E_{11} = A^{-1} + A^{-1} B E_{22} C A^{-1}$
- $E_{12} = -A^{-1} B E_{22}$
- $E_{21} = -E_{22} C A^{-1}$
- $E_{22} = (D - C A^{-1} B)^{-1}$



# 中心化

- 中心化的本质

- 响应变量

$$\mathbf{y}^* = \mathbf{y} - \mathbf{1}_n(\mathbf{1}_n^\top \mathbf{1}_n)^{-1} \mathbf{1}_n^\top \mathbf{y} = (\mathbf{I}_n - \mathbf{H}_{\mathbf{1}_n})\mathbf{y}$$

- 自变量

$$\mathbf{X}_c = \mathbf{X}_o - \mathbf{1}_n(\mathbf{1}_n^\top \mathbf{1}_n)^{-1} \mathbf{1}_n^\top \mathbf{X}_o = (\mathbf{I}_n - \mathbf{H}_{\mathbf{1}_n})\mathbf{X}_o$$

- 提问：为什么  $\mathbf{H}_{\mathbf{1}_n} = \mathbf{1}_n(\mathbf{1}_n^\top \mathbf{1}_n)^{-1} \mathbf{1}_n^\top$  是一个对称幂等矩阵？

# 中心化

- 估计  $\beta_0$  的影响

$$\begin{aligned}\hat{\beta}_{c,0} &= (n^{-1}\mathbf{1}_n^\top + n^{-2}\mathbf{1}_n^\top \mathbf{X}_c \mathbf{A}_c \mathbf{X}_c^\top \mathbf{1}_n \mathbf{1}_n^\top - n^{-1}\mathbf{1}_n^\top \mathbf{X}_c \mathbf{A}_c \mathbf{X}_c^\top) \mathbf{y}^* \\ &= (n^{-1}\mathbf{1}_n^\top + n^{-2}\mathbf{1}_n^\top (\mathbf{I}_n - \mathbf{H}_{\mathbf{1}_n}) \mathbf{X}_o \mathbf{A}_c \mathbf{X}_c^\top \mathbf{1}_n \mathbf{1}_n^\top - n^{-1}\mathbf{1}_n^\top (\mathbf{I}_n - \mathbf{H}_{\mathbf{1}_n}) \mathbf{X}_o \mathbf{A}_c \mathbf{X}_c^\top) \mathbf{y}^* \\ &= n^{-1}\mathbf{1}_n^\top \mathbf{y}^* \\ &= n^{-1}\mathbf{1}_n^\top (\mathbf{I}_n - \mathbf{H}_{\mathbf{1}_n}) \mathbf{y} \\ &= 0\end{aligned}$$

计算提示

$$\mathbf{1}_n^\top (\mathbf{I}_n - \mathbf{H}_{\mathbf{1}_n}) = \mathbf{1}_n^\top - \mathbf{1}_n^\top \mathbf{H}_{\mathbf{1}_n} = \mathbf{1}_n^\top - \mathbf{1}_n^\top \mathbf{1}_n (\mathbf{1}_n^\top \mathbf{1}_n)^{-1} \mathbf{1}_n^\top = \mathbf{1}_n^\top - \mathbf{1}_n^\top = \mathbf{0}$$

# 中心化

- 估计  $\beta_1, \dots, \beta_p$  的影响

$$\begin{aligned}\hat{\beta}_{c,1} &= (-n^{-1}A_cX_c^\top \mathbf{1}_n \mathbf{1}_n^\top + A_cX_c^\top) \mathbf{y}^* \\ &= (-n^{-1}A_cX_o^\top (\mathbf{I}_n - H_{\mathbf{1}_n}) \mathbf{1}_n \mathbf{1}_n^\top + A_cX_c^\top) \mathbf{y}^* \\ &= A_cX_c^\top \mathbf{y}^* \\ &= A_oX_o^\top (\mathbf{I}_n - H_{\mathbf{1}_n}) \mathbf{y}\end{aligned}$$

$$\begin{aligned}\hat{\beta}_1 &= (-n^{-1}A_oX_o^\top \mathbf{1}_n \mathbf{1}_n^\top + A_oX_o^\top) \mathbf{y} \\ &= A_oX_o^\top (\mathbf{I}_n - \mathbf{1}_n (\mathbf{1}_n^\top \mathbf{1}_n)^{-1} \mathbf{1}_n^\top) \mathbf{y} \\ &= A_oX_o^\top (\mathbf{I}_n - H_{\mathbf{1}_n}) \mathbf{y}\end{aligned}$$

## 计算提示

$$\begin{aligned}A_c &= (X_c^\top X_c - n^{-1}X_c^\top \mathbf{1}_n \mathbf{1}_n^\top X_c)^{-1} = (X_o^\top (\mathbf{I}_n - H_{\mathbf{1}_n}) X_o - n^{-1}X_o^\top (\mathbf{I}_n - H_{\mathbf{1}_n}) \mathbf{1}_n \mathbf{1}_n^\top (\mathbf{I}_n - H_{\mathbf{1}_n}) X_o)^{-1} \\ &= (X_o^\top (\mathbf{I}_n - H_{\mathbf{1}_n}) X_o)^{-1} = A_o\end{aligned}$$

小结

		原始数据	中心化	标准化
数据变换	响应变量	$y = (y_1, \cdots, y_n)^\top$	$\mathbf{y}^* = (y_1^*, \cdots, y_n^*)^\top$ $y_i^* = y_i - \bar{y}, \bar{y} = n^{-1} \sum_{i=1}^n y_i$	$\mathbf{y}^{**} = (y_1^{**}, \cdots, y_n^{**})^\top$ $y_i^{**} = \frac{y_i^*}{\sqrt{l_{yy}}}, l_{yy} = \sum_{i=1}^n (y_i^*)^2$
	自变量	$\mathbf{X} = (\mathbf{1}_n, \mathbf{X}_o),$ $\mathbf{X}_o = \{x_{ij}\}_{n \times p}$	$\mathbf{X}^* = (\mathbf{1}_n, \mathbf{X}_c), \mathbf{X}_c = \{x_{ij}^*\}_{n \times p}$ $x_{ij}^* = x_{ij} - \bar{x}_j, \bar{x}_j = n^{-1} \sum_{i=1}^n x_{ij}$	$\mathbf{X}^{**} = (\mathbf{1}_n, \mathbf{X}_s), \mathbf{X}_s = \{x_{ij}^{**}\}_{n \times p}$ $x_{ij}^{**} = \frac{x_{ij}^*}{\sqrt{l_{jj}}}, l_{jj} = \sum_{i=1}^n (x_{ij}^*)^2$
参数估计	$\beta_0$	$\hat{\beta}_0$	$\hat{\beta}_{c,0} = 0$	$\hat{\beta}_{s,0} = ?$
	$\beta_1$	$\hat{\beta}_1$	$\hat{\beta}_{c,1} = \hat{\beta}_1$	$\hat{\beta}_{s,1} = ?$

# 标准化

- 定义

$$L = \begin{pmatrix} \frac{1}{\sqrt{l_{11}}} & & \\ & \ddots & \\ & & \frac{1}{\sqrt{l_{pp}}} \end{pmatrix}$$

- 于是有

$$\boldsymbol{X}_S = \boldsymbol{X}_c \boldsymbol{L}$$

$$\boldsymbol{y}^{**} = \frac{1}{\sqrt{l_{yy}}} \boldsymbol{y}^*$$

# 标准化

• 估计  $\beta_1 = (\beta_1, \dots, \beta_p)^\top$  的影响

$$\begin{aligned}
 \hat{\beta}_{s,1} &= (X_s^\top X_s)^{-1} X_s^\top y^{**} \\
 &= (LX_c^\top X_c L)^{-1} L X_c^\top \left( \frac{1}{\sqrt{l_{yy}}} y^* \right) \\
 &= L^{-1} (X_c^\top X_c)^{-1} L^{-1} L X_c^\top \left( \frac{1}{\sqrt{l_{yy}}} y^* \right) \\
 &= \frac{1}{\sqrt{l_{yy}}} L^{-1} (X_c^\top X_c)^{-1} X_c^\top y^*
 \end{aligned}$$

$$\begin{aligned}
 \hat{\beta}_{c,1} &= (-n^{-1} A_c X_c^\top \mathbf{1}_n \mathbf{1}_n^\top + A_c X_c^\top) y^* \\
 &= A_c X_c^\top y^* \\
 &= (X_c^\top X_c - n^{-1} X_c^\top \mathbf{1}_n \mathbf{1}_n^\top X_c)^{-1} X_c^\top y^* \\
 &= (X_c^\top X_c)^{-1} X_c^\top y^*
 \end{aligned}$$

# 标准化

- 估计  $\beta_1 = (\beta_1, \dots, \beta_p)^\top$  的影响

$$\hat{\beta}_{s,1} = \frac{1}{\sqrt{l_{yy}}} L^{-1} \hat{\beta}_{c,1} = \frac{1}{\sqrt{l_{yy}}} L^{-1} \hat{\beta}_1$$

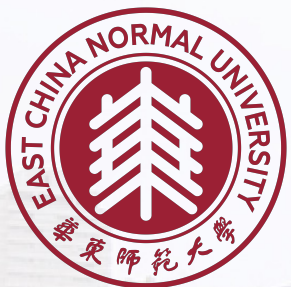
其分量为

$$\hat{\beta}_{sj} = \frac{\sqrt{l_{jj}}}{\sqrt{l_{yy}}} \hat{\beta}_{cj} = \frac{\sqrt{l_{jj}}}{\sqrt{l_{yy}}} \hat{\beta}_j, j = 1, 2, \dots, p$$

# 总结

		原始数据	中心化	标准化
数据变换	响应变量	$\mathbf{y} = (y_1, \cdots, y_n)^\top$	$\mathbf{y}^* = (y_1^*, \cdots, y_n^*)^\top$ $y_i^* = y_i - \bar{y}, \bar{y} = n^{-1} \sum_{i=1}^n y_i$	$\mathbf{y}^{**} = (y_1^{**}, \cdots, y_n^{**})^\top$ $y_i^{**} = \frac{y_i^*}{\sqrt{l_{yy}}}, l_{yy} = \sum_{i=1}^n (y_i^*)^2$
	自变量	$\mathbf{X} = (\mathbf{1}_n, \mathbf{X}_o),$ $\mathbf{X}_o = \{x_{ij}\}_{n \times p}$	$\mathbf{X}^* = (\mathbf{1}_n, \mathbf{X}_c), \mathbf{X}_c = \{x_{ij}^*\}_{n \times p}$ $x_{ij}^* = x_{ij} - \bar{x}_j, \bar{x}_j = n^{-1} \sum_{i=1}^n x_{ij}$	$\mathbf{X}^{**} = (\mathbf{1}_n, \mathbf{X}_s), \mathbf{X}_s = \{x_{ij}^{**}\}_{n \times p}$ $x_{ij}^{**} = \frac{x_{ij}^*}{\sqrt{l_{jj}}}, l_{jj} = \sum_{i=1}^n (x_{ij}^*)^2$
参数估计	$\beta_0$	$\hat{\beta}_0$	$\hat{\beta}_{c,0} = 0$	$\hat{\beta}_{s,0} = 0$
	$\begin{pmatrix} \beta_1 \\ \vdots \\ \beta_p \end{pmatrix}$	$\hat{\boldsymbol{\beta}}_1$	$\hat{\boldsymbol{\beta}}_{c,1} = \hat{\boldsymbol{\beta}}_1$	$\hat{\boldsymbol{\beta}}_{s,1} = \frac{1}{\sqrt{l_{yy}}} \begin{pmatrix} \sqrt{l_{11}} & & \\ & \ddots & \\ & & \sqrt{l_{pp}} \end{pmatrix} \hat{\boldsymbol{\beta}}_1$





# 谢谢



SCHOOL OF DATA  
SCIENCE & ENGINEERING  
数据科学与工程学院