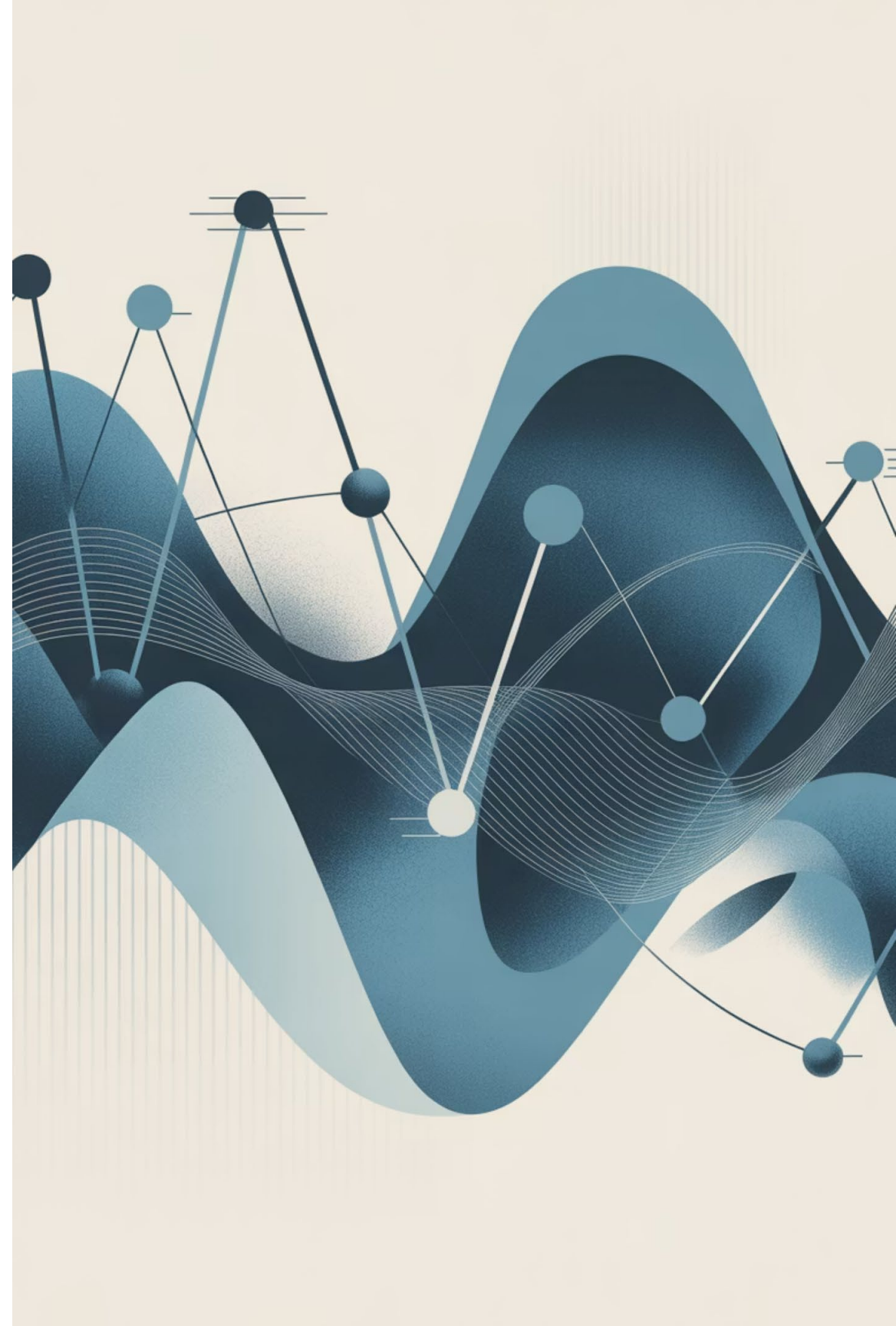


# 相关分析与多元 线性回归

统计建模的基本原理与应用



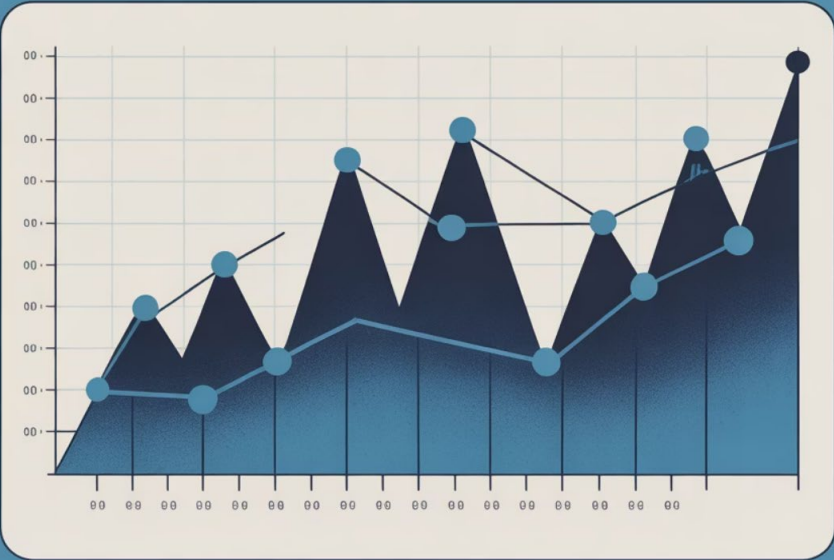
# 相关分析:揭示变量间的线性关系

相关分析用于研究两个或多个变量之间**线性关系**的强度与方向,回答"是否有关""有多强、方向如何"等核心问题。通过计算皮尔逊相关系数,我们可以快速识别变量间的关联模式。

相关系数的数学表达式为:

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

其中 $r \in [-1, 1]$ ,取值范围反映了关系的强度和方向。



# 相关系数的解读标准

强相关

$|r| \geq 0.8$

变量间存在显著的线性关系,一个变量的变化能很好地预测另一个变量

中等相关

$0.5 \leq |r| < 0.8$

变量间有明显关联,但存在一定程度的波动和不确定性

弱相关

$0.3 \leq |r| < 0.5$

变量间关系较弱,可能需要考虑其他影响因素

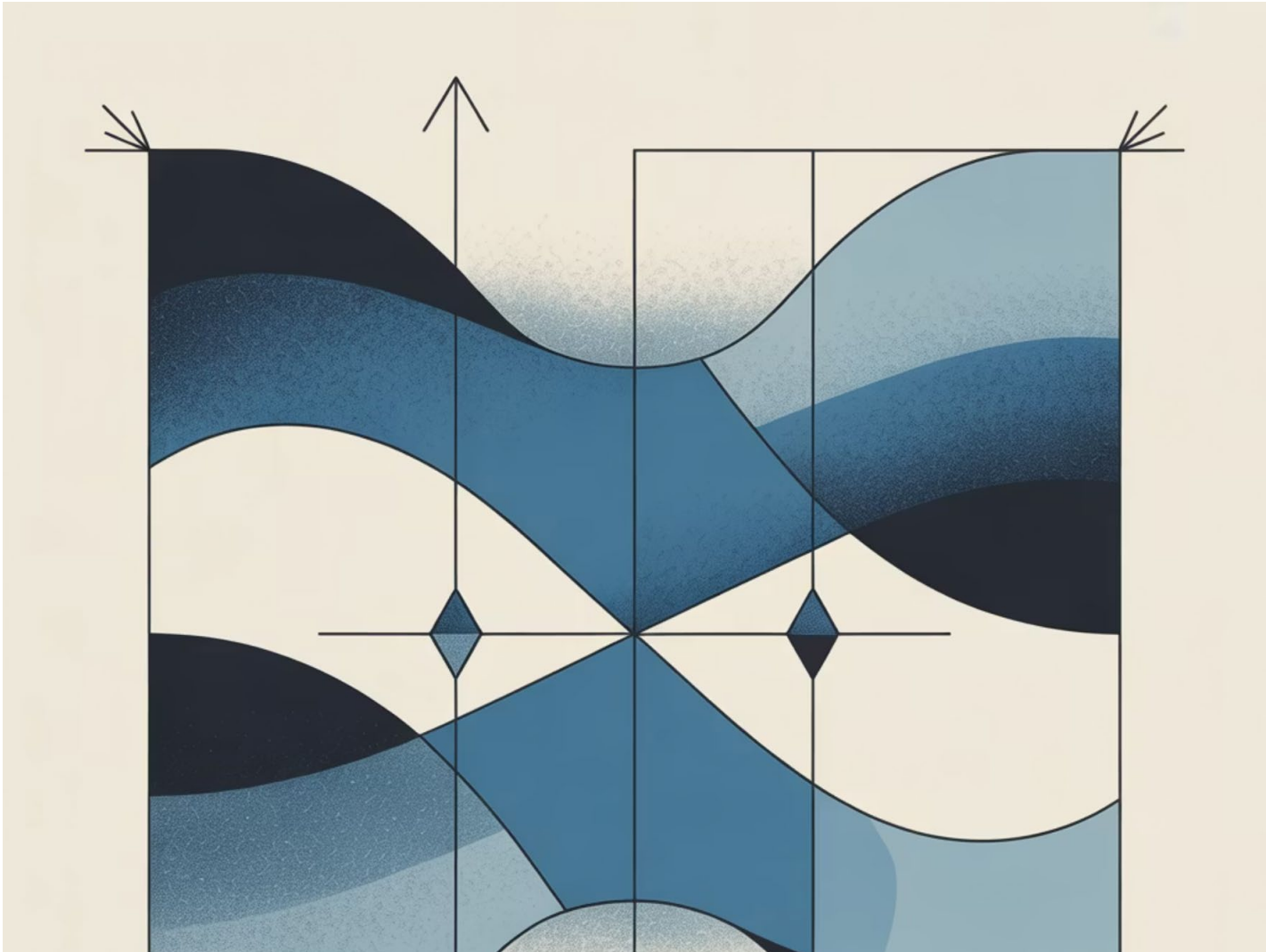
几乎无关

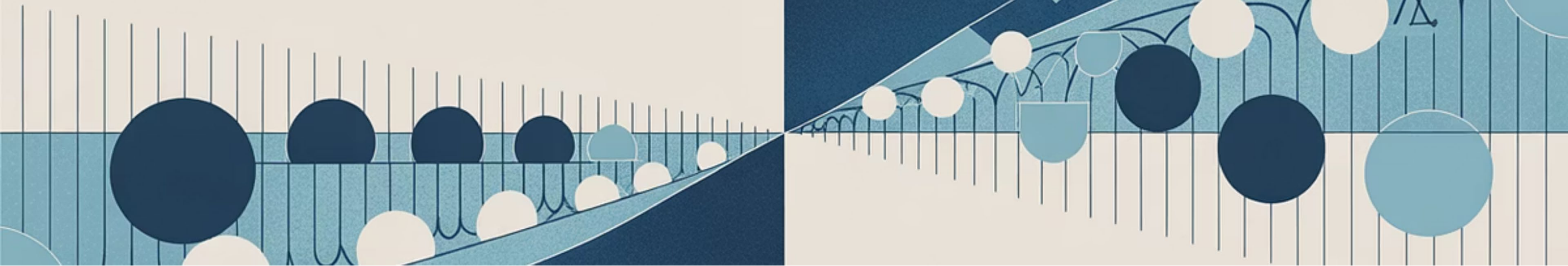
$|r| < 0.3$

变量间线性关系极弱或不存在,需探索非线性关系

## 方向判断

- 正相关 ( $r > 0$ ): 变量同向变化
- 负相关 ( $r < 0$ ): 变量反向变化
- 无相关 ( $r = 0$ ): 无线性关系





# 相关分析的优势与局限

## 核心优势

- 简单直观,易于理解和解释
- 快速识别变量间的关系模式
- 为后续建模提供明确方向
- 计算效率高,适合大规模筛选

## 重要局限

- 仅能识别线性关系,无法捕捉非线性模式
- 相关不等于因果,不能推断变量间的因果关系
- 对异常值高度敏感,可能产生误导
- 无法量化单个变量的边际影响

# 多元线性回归:量化变量关系

多元线性回归在相关分析基础上更进一步,不仅识别关系的存在,更**量化**因变量与多个自变量之间的线性关系,为预测和决策提供精确的数学模型。



## 发现关系

相关分析识别变量间是否存在关联



## 量化影响

回归分析计算每个变量的具体影响程度



## 预测应用

建立预测模型,支持实际决策

回归模型的数学表达式:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_k X_k + \varepsilon$$

其中  $Y$  为因变量,  $X_1, \dots, X_k$  为自变量,  $\beta_0$  为截距,  $\beta_1, \dots, \beta_k$  为回归系数,  $\varepsilon$  为随机误差项。



# 最小二乘法:参数估计的核心

## 估计原理

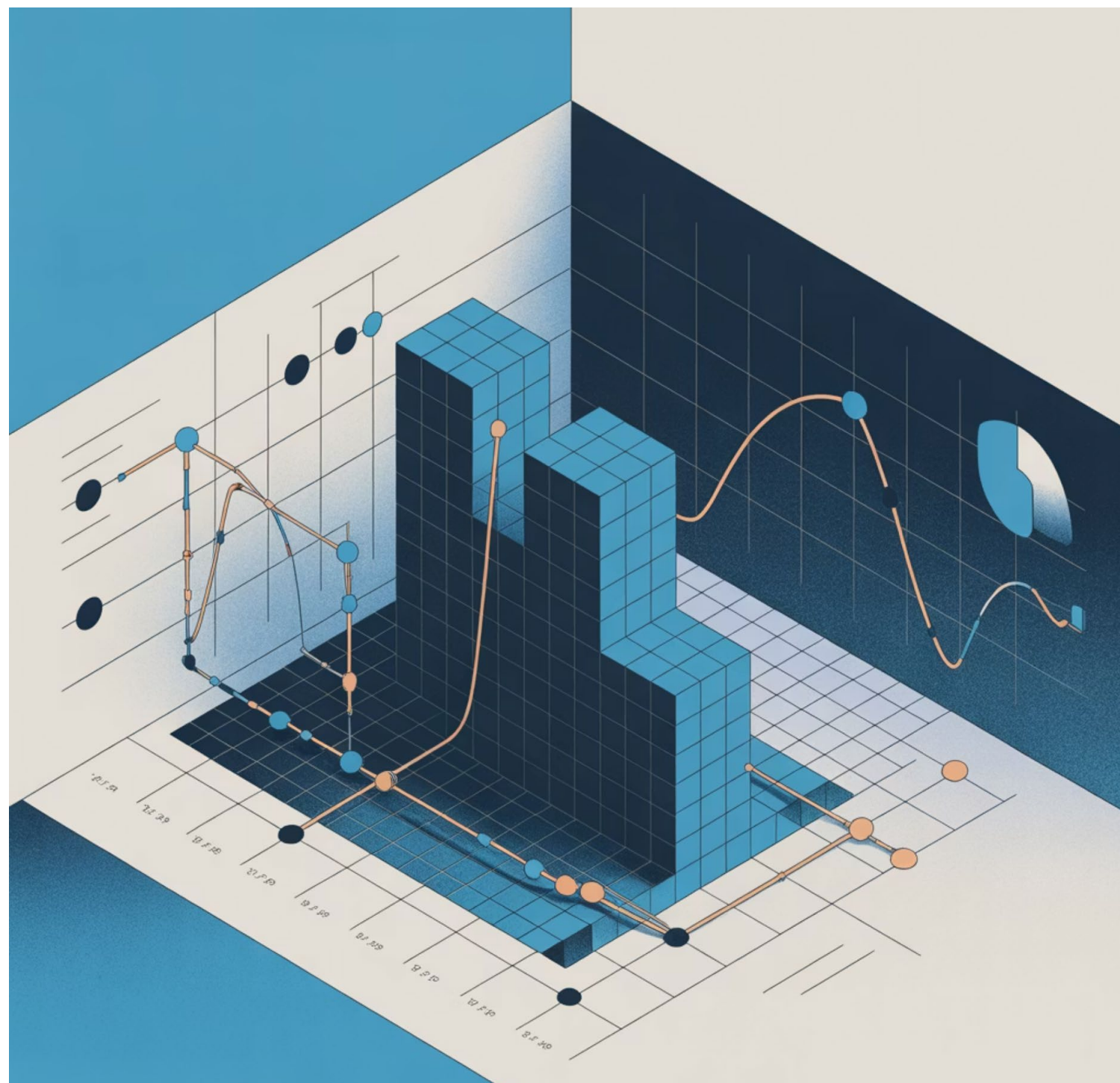
最小二乘法(OLS)通过最小化观测值与预测值之间的平方误差和来估计回归参数,确保模型对数据的最佳拟合。

目标函数:

$$\min_{\beta_0, \dots, \beta_k} \sum_{i=1}^n [y_i - (\beta_0 + \beta_1 x_{1i} + \dots + \beta_k x_{ki})]^2$$

矩阵解:

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



# 经典线性回归模型的五大假设

( )

## 线性假设

因变量  $Y$  与各自变量  $X_j$  之间的关系可以用线性方程表示,模型形式正确

( [

## 独立性假设

各观测值之间相互独立,不存在自相关或序列相关问题

( {

## 同方差假设

误差项的方差保持恒定,即  $Var(\varepsilon_i) = \sigma^2$

( }

## 正态性假设

误差项服从正态分布  $\varepsilon_i \sim N(0, \sigma^2)$ ,用于小样本统计推断

( @

## 无多重共线性

自变量之间不存在完全线性依赖关系,保证参数估计的稳定性

# 模型评估:拟合优度与显著性检验



## 决定系数 $R^2$

$$R^2 = 1 - \frac{SSE}{SST} = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

取值范围  $R^2 \in [0, 1]$ , 越接近1表示模型拟合效果越好



## 调整决定系数

$$\bar{R}^2 = 1 - (1 - R^2) \frac{n - 1}{n - k - 1}$$

考虑自变量数量的影响,避免过度拟合,更适合比较不同模型

## 整体显著性检验 (F检验)

$$F = \frac{SSR/k}{SSE/(n - k - 1)}$$

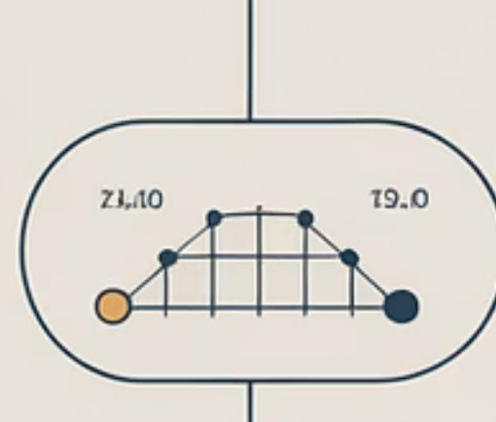
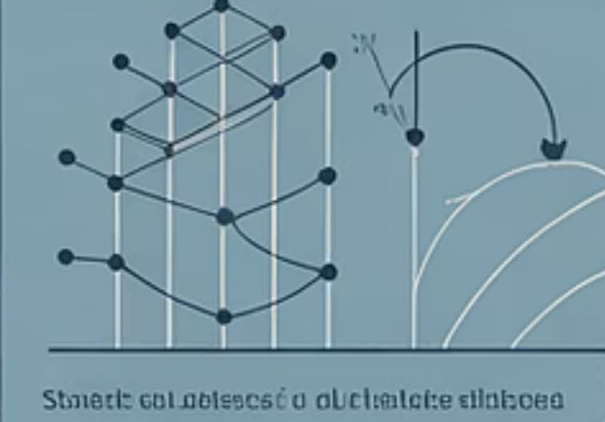
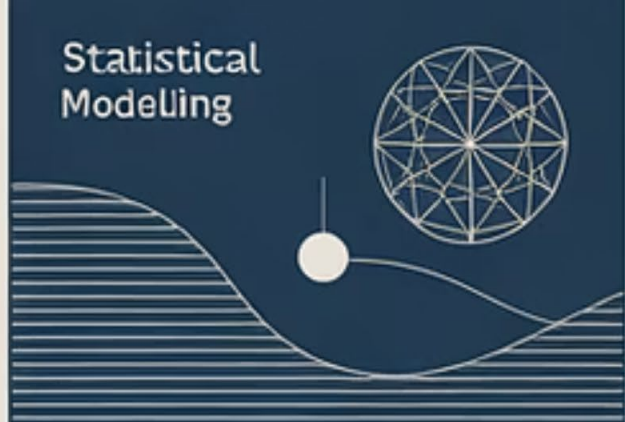
检验所有自变量联合对因变量的解释能力是否显著

## 个别显著性检验 (t检验)

$$t = \frac{\hat{\beta}_i}{SE(\hat{\beta}_i)}$$

检验单个自变量对因变量的影响是否显著





# 实践应用:变量选择与模型诊断

## Step 1: 相关性筛选

计算各变量与因变量的相关系数,初步识别潜在预测变量

1

## Step 3: 共线性检查

检查自变量间相关性,处理多重共线性问题

3

## Step 2: 低相关剔除

剔除相关度过低的变量(如 $|r| < 0.3$ ),提高模型效率

2

## Step 4: 建模优化

进入回归建模阶段,进行迭代优化和验证

4

## 模型诊断的三大支柱



### 残差分析

通过残差-拟合值图检验同方差性,Q-Q图检验正态性,残差-自变量图验证线性假设



### 共线性诊断

使用VIF(方差膨胀因子)、容忍度和条件指数识别自变量间的多重共线性



### 独立性检验

采用Durbin-Watson统计量检验误差项的自相关性,确保观测独立

# 相关分析与回归分析的协同作用



**核心要点:**相关分析与回归分析形成递进关系——从描述关系到量化关系,再到预测应用。两者互补协作,相关分析用于发现和筛选,回归分析用于量化和预测,共同构成完整的统计建模体系。