

# 解决多重共线性的代码实现

检测方法以及解决方法

小胖



# 目录

## ONE 两个变量的检测方法

one-way ANOVA

## TWO 多变量的检测方法

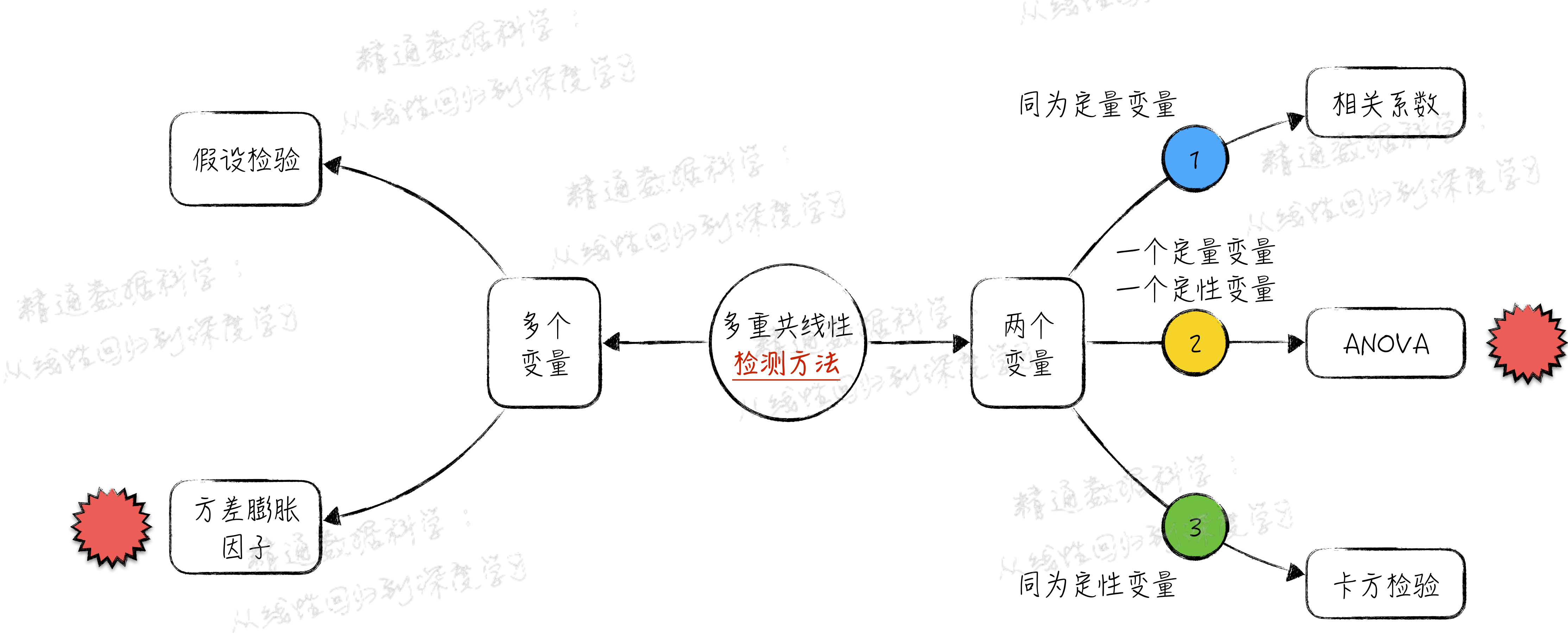
方差膨胀因子

## THREE 虚拟变量陷阱

一个简单的例子

# 两个变量的检测方法

整体回顾

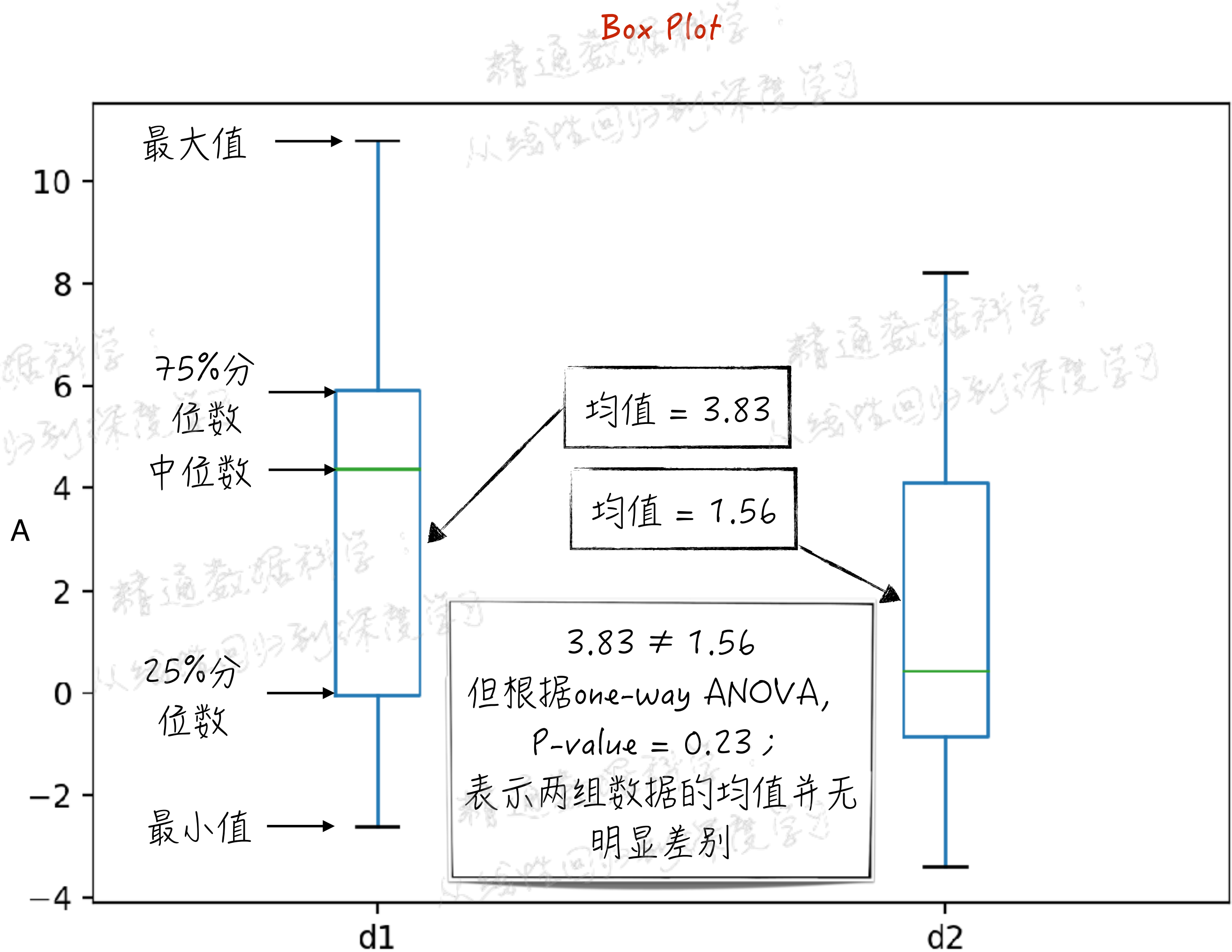


# 两个变量的检测方法

one-way ANOVA

one-way ANOVA: 多组数据的均值是否显著不同

ANOVA : analysis of variance



d1, d2同为均值等于5的正态分布



# 两个变量的检测方法

one-way ANOVA

ANOVA : analysis of variance

多重共线性的检测

假设A是定量变量、B为定性变量（共有k个分类）

变量A按B的取值分为k组

定义技术指标  $\eta^2$

变量A的**全部**方差  $S_{total}$

变量A的**组内**方差  $S_{within}$

变量A的**组间**方差  $S_{between}$

$$\eta^2 = \frac{S_{between}}{S_{total}}$$

$\eta^2$  越接近1，共线性越严重

# 目录

## ONE 两个变量的检测方法

one-way ANOVA

## TWO 多变量的检测方法

方差膨胀因子

## THREE 虚拟变量陷阱

一个简单的例子

# 多变量的检测方法

方差膨胀因子

假设如下的线性回归模型： $y = \beta_0 + \beta_1x_1 + \dots + \beta_kx_k + \varepsilon$

随机扰动项  
的方差

数学上可以证明：

$$\hat{Var}(\hat{\beta}_1) = \frac{\sigma^2}{Var(x_1)} \frac{1}{1 - R_1^2}$$

下面线性回  
归模型的决  
定系数

$x_1$ 能多大程度被其他变量解释： $x_1 = \alpha_0 + \alpha_2x_2 + \dots + \alpha_kx_k + e$

计算VIF时，不要忘记在  
数据里加入常数项

定义方差膨胀因子：

$$VIF_1 = \frac{1}{1 - R_1^2}$$

只与数据相关，  
与使用的模型无关

# 目录

## ONE 两个变量的检测方法

one-way ANOVA

## TWO 多变量的检测方法

方差膨胀因子

## THREE 虚拟变量陷阱

一个简单的例子



# 虚拟变量陷阱

这里是副标题文字

定性变量有四个取值：a, b, c, d

定义各类别个数

选择a作为基准类别，检测共线性

选择b作为基准类别，检测共线性

类别	个数
a	5
b	100
c	70
d	20

# THANK YOU

精通数据挖掘科学：  
从线性回归到深度学习