

# 第16章 不平衡数据问题及处理

《Python数据科学：技术详解与商业实践》

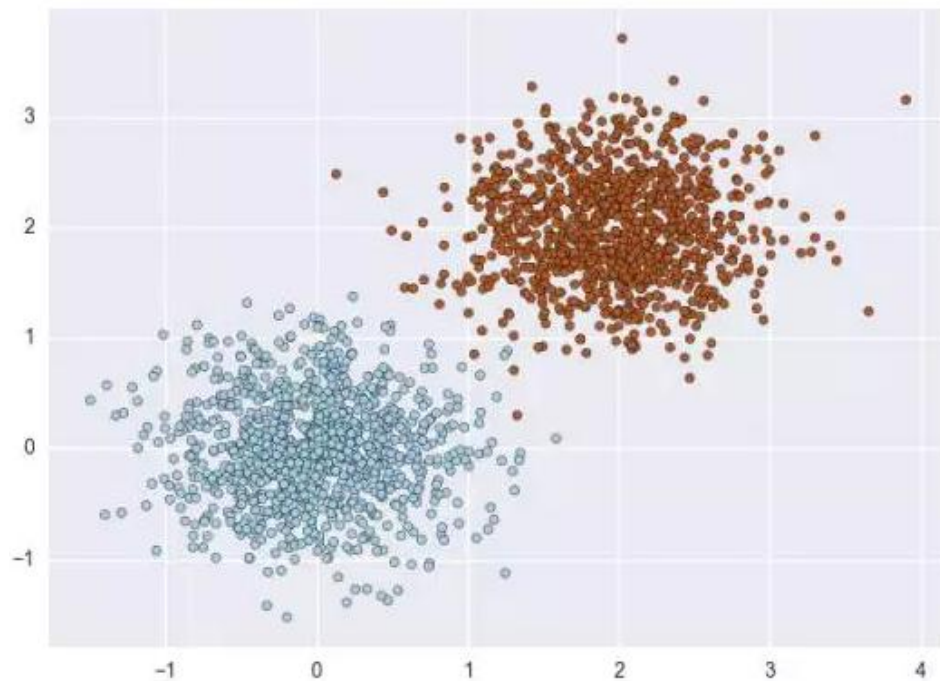
讲师：Ben

# 自我介绍

- 天善商业智能和大数据社区      讲师 – Ben
- 天善社区 ID - Ben\_Chang
- <https://www.hellobi.com> – 学习过程中有任何相关的问题都可以提到技术社区数据挖掘版块。

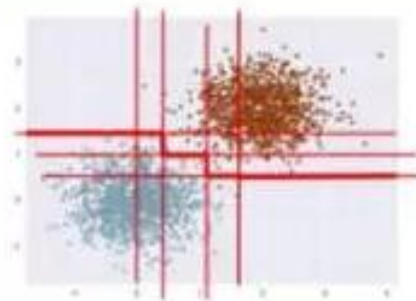
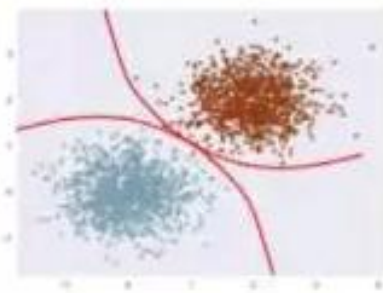
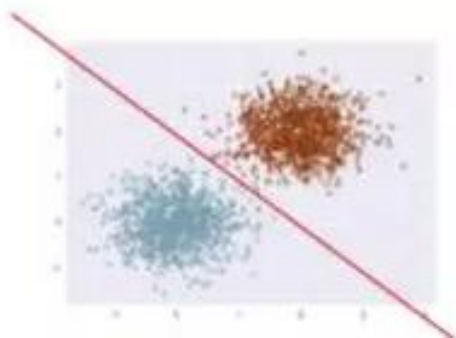
- 不平衡分类概述
- 欠采样法
- 过采样法
- 综合采样

## 以往学习的算法：



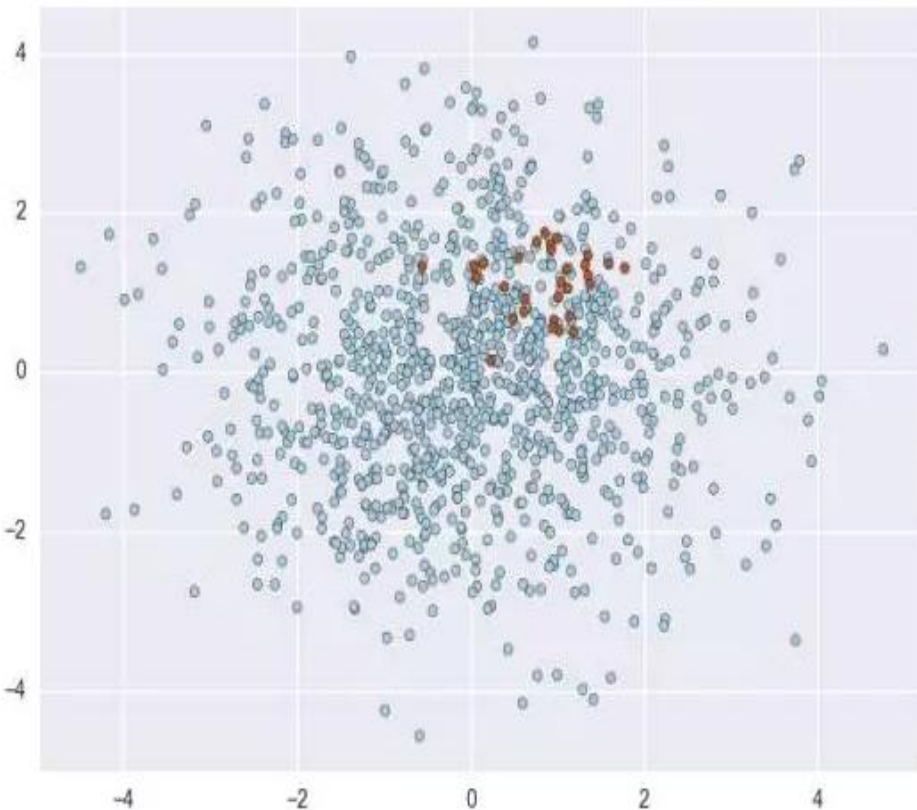
在以上章节中均认为数据是对称分布的，即正负样本的数量相当。这样可以把注意力集中特定算法上，而不被其它问题干扰。

分类算法的目标是尝试学习出一个能够分辨二者的分离器（分类器）。根据不同的数学、统计或几何假设，达成这一目标的方法很多：



# 不平衡分类数据

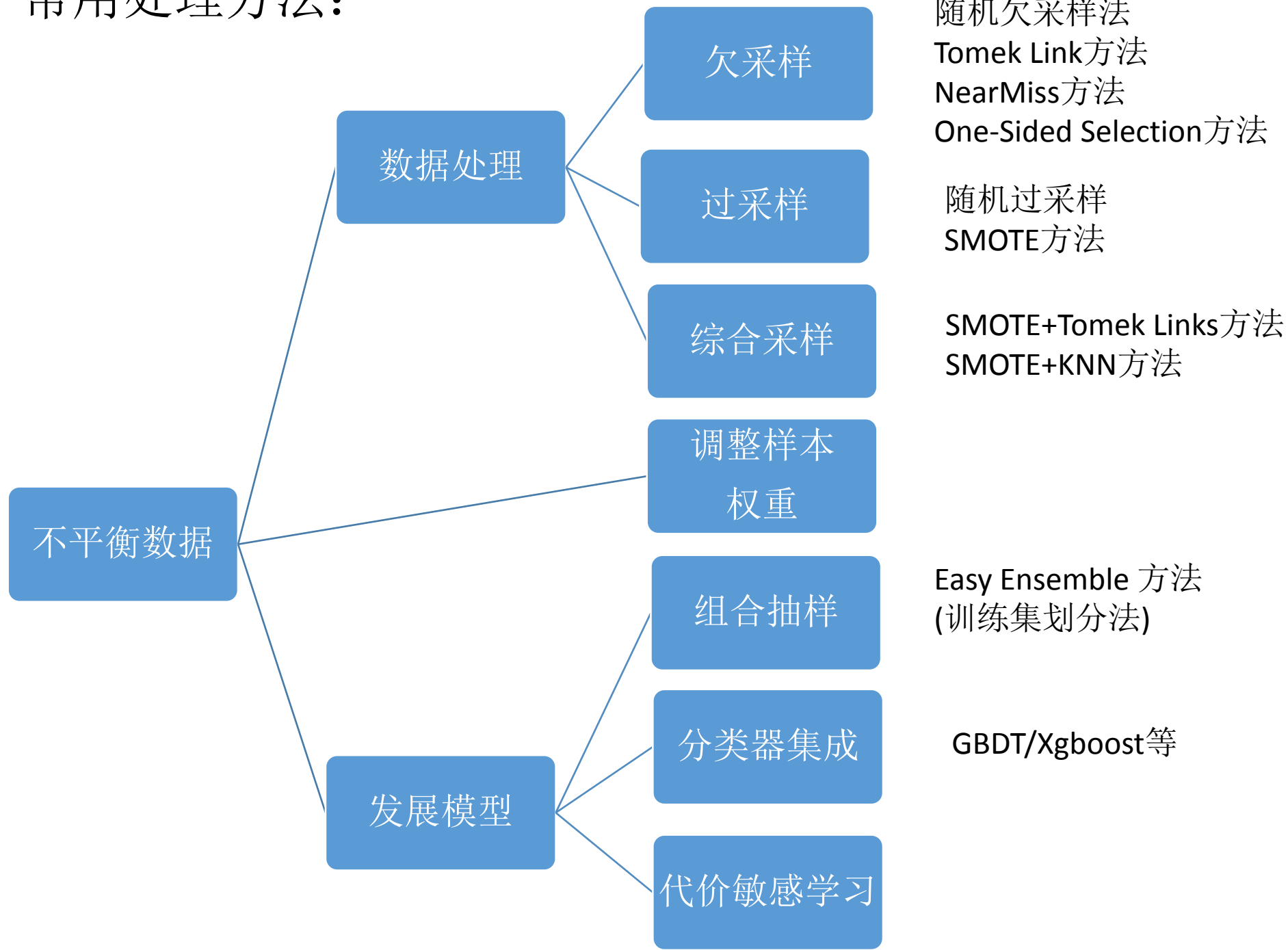
然而，当你开始面对真实的、未加工过的数据时，你会马上注意到，这些数据要嘈杂且不平衡得多。真实数据的散点图看起来更像是这样的：



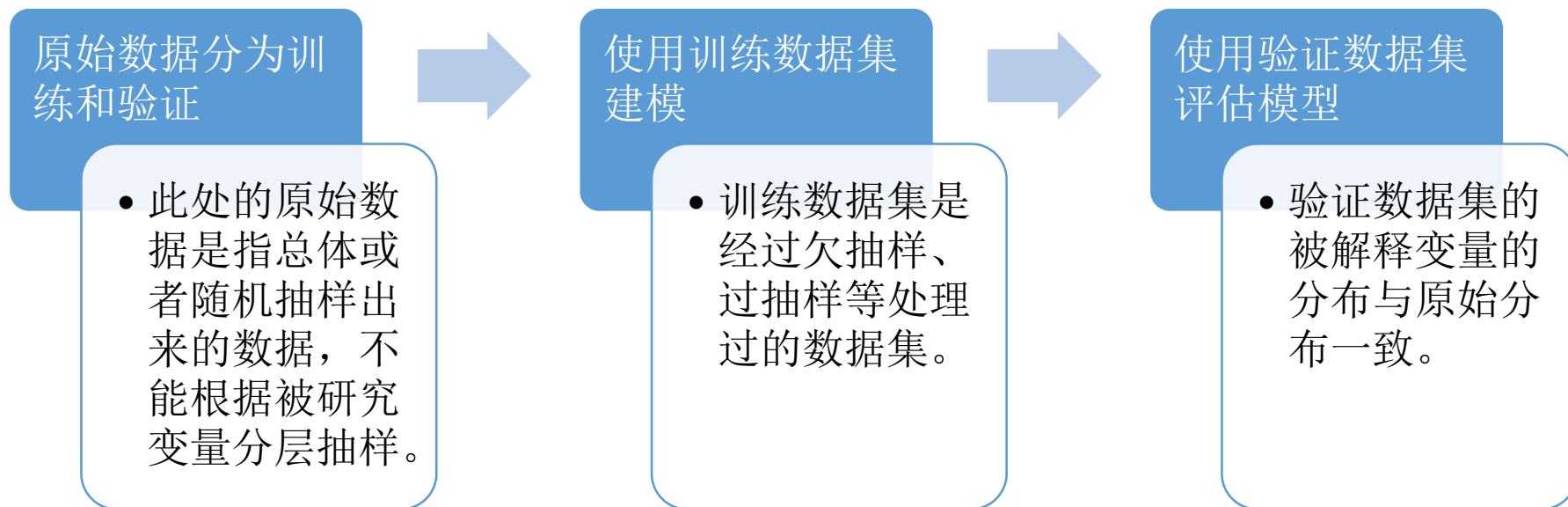
对于不平衡类的研究通常认为「不平衡」意味着少数类只占 **10% 到 20%**。而在现实中，数据库甚至能够比上面的例子更加不平衡。以下是一些例子：

- 每年，约 **2%** 的信用卡账户是伪造的（多数的欺诈检测领域是极其不平衡的）；
- 针对某一病征的医学筛查通常涵盖了许多没有此病征的人，以检查出少数患者（例：美国的 HIV 感染率约为 **0.4%**）；
- 每年，硬盘驱动器故障的发生率约为 **1%**；
- 在线广告的转化率在  $10^{-3}$  到  $10^{-6}$  的范围内；
- 工厂的产品缺陷率一般在 **0.1%** 左右。

# 常用处理方法:



# 数据处理方法的全流程：

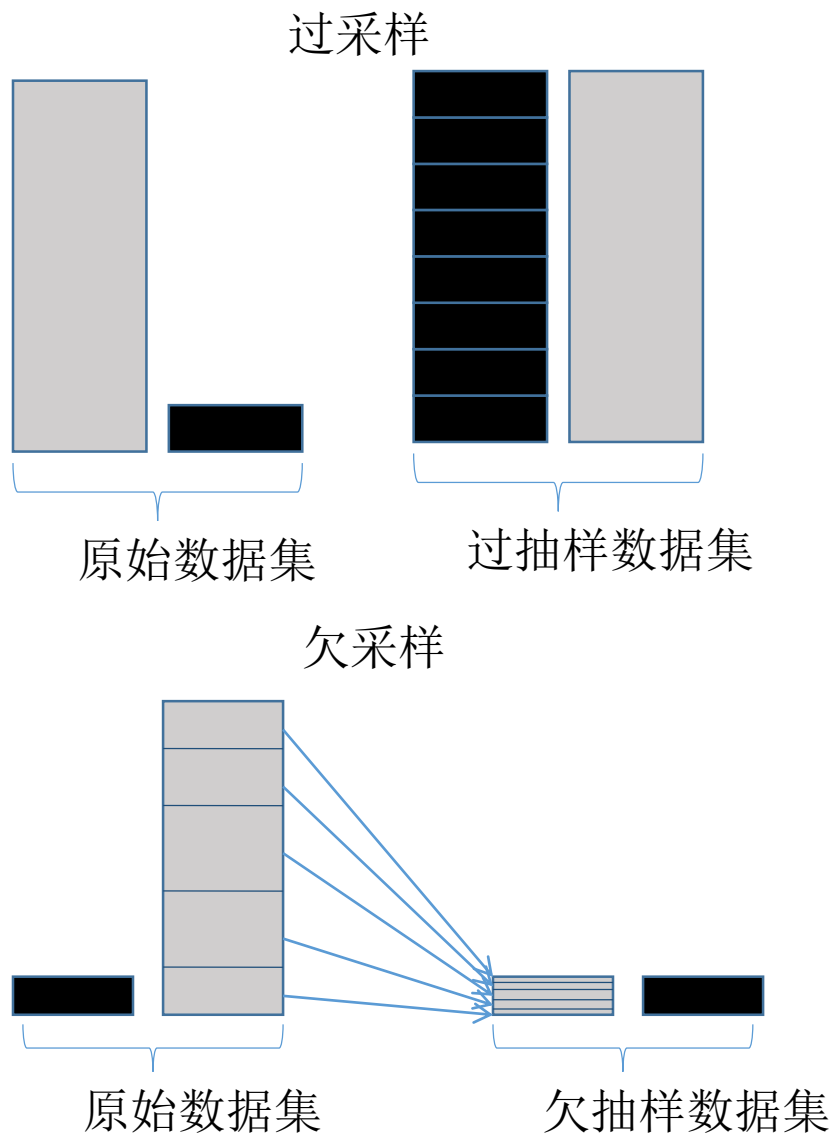


## 注意事项：

- 1、评估指标：使用精确度（Precise Rate）、召回率（Recall Rate）、Fmeasure或ROC曲线、准确度召回曲线（precision-recall curve）；不要使用准确度（Accurate Rate）
- 2、不要使用模型给出的标签，而是要概率估计；得到概率估计之后，不要盲目地使用0.50的决策阈值来区分类别，应该再检查表现曲线之后再自己决定使用哪个阈值。



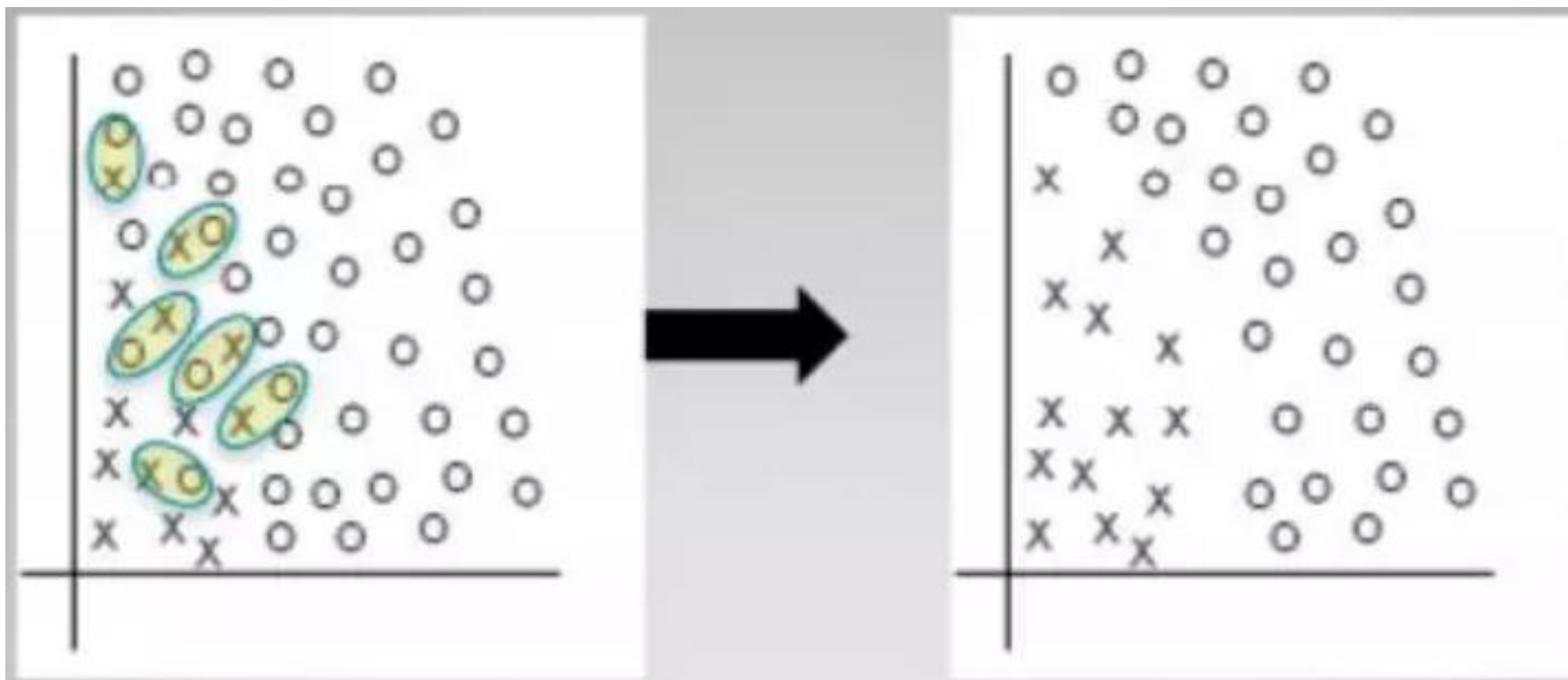
# 随机过抽样与欠抽样



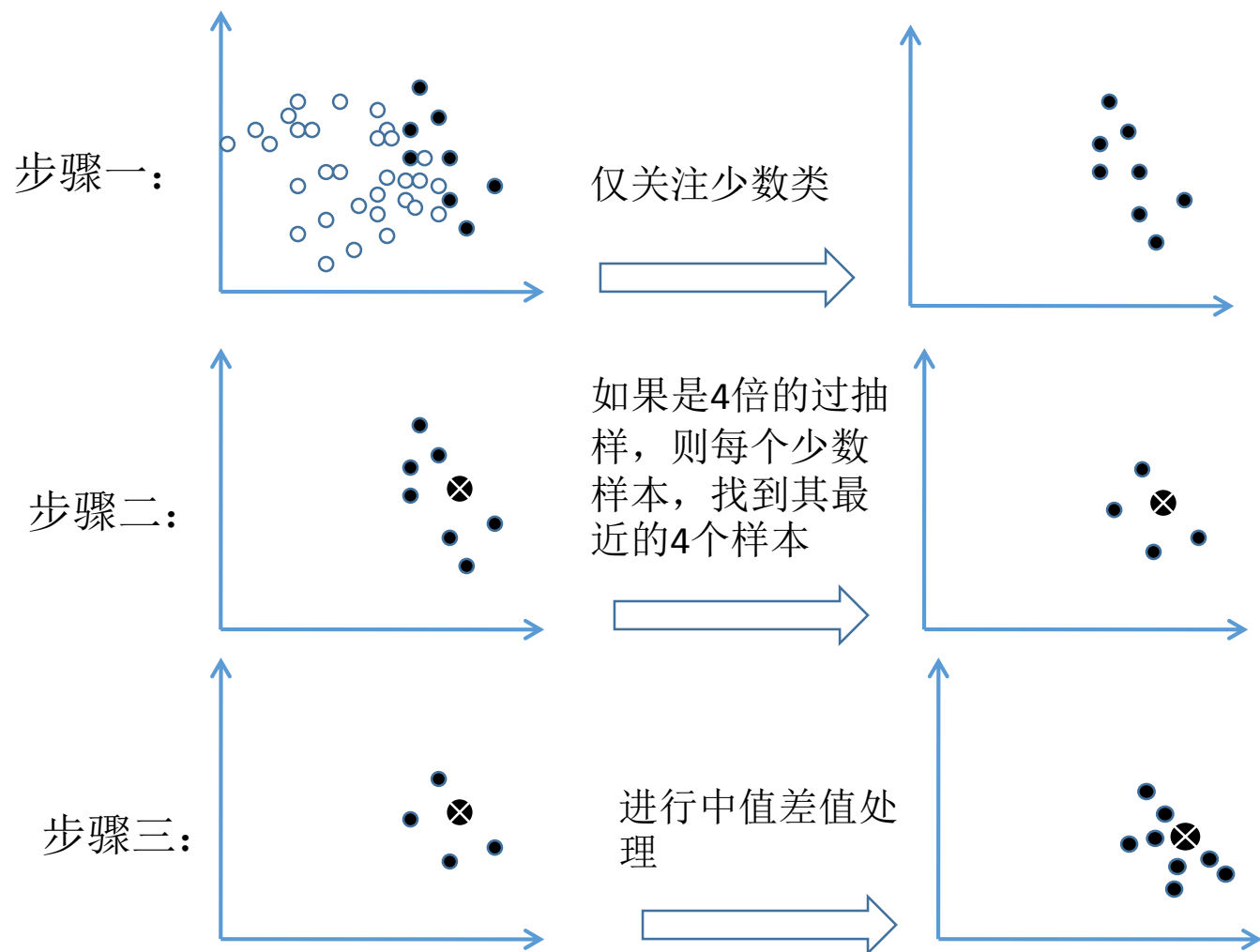
过采样会随机复制少数样例以增大它们的规模。欠采样则随机地少采样主要的类。一些数据科学家（天真地）认为过采样更好，因为其会得到更多的数据，而欠采样会将数据丢掉。但请记住复制数据不是没有后果的——因为其会得到复制出来的数据，它就会使变量的方差表面上比实际上更小。而过采样的好处是它也会复制误差的数量：如果一个分类器在原始的少数类数据集上做出了一个错误的负面错误，那么将该数据集复制五次之后，该分类器就会在新的数据集上出现六个错误。相对地，欠采样会让独立变量（independent variable）的方差看起来比其实际的方差更高。



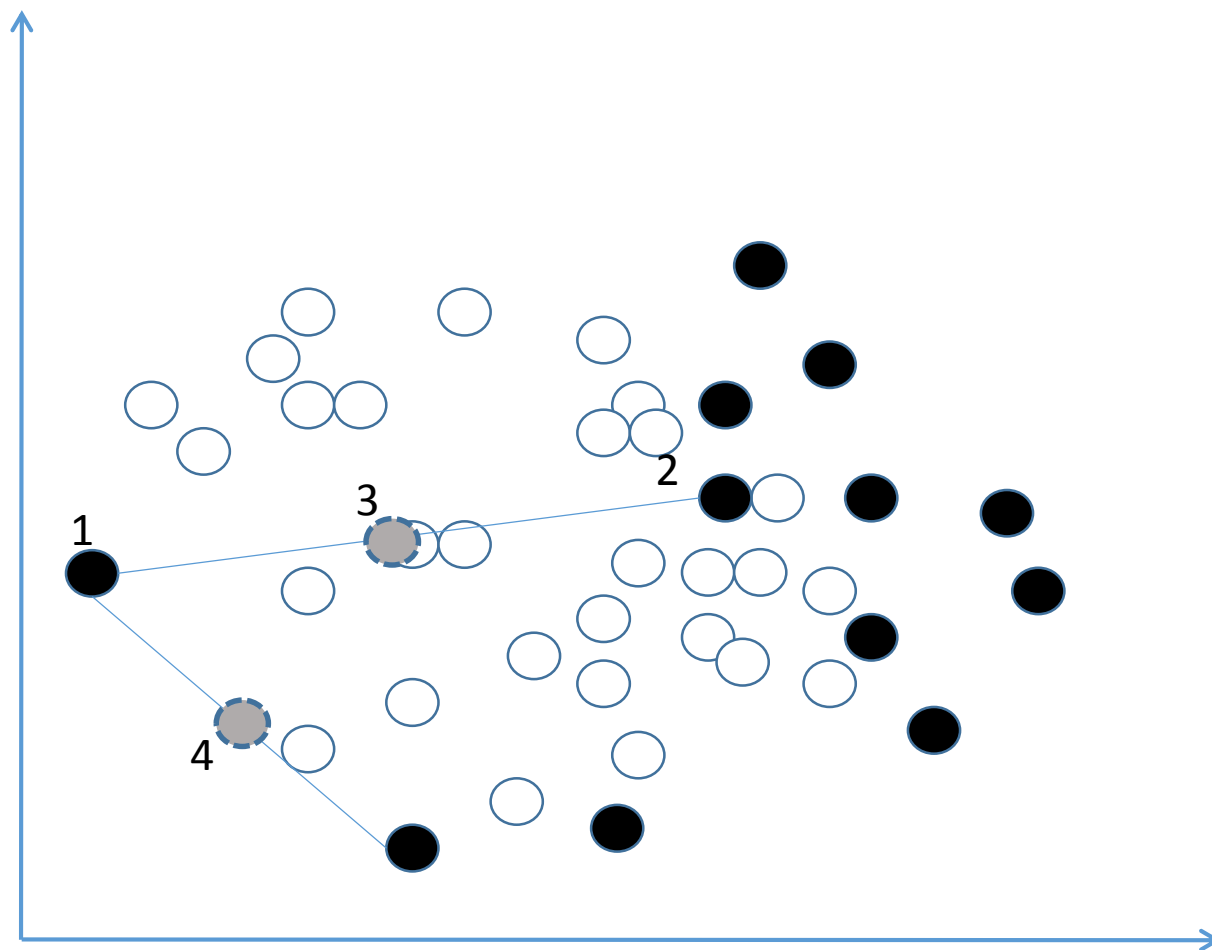
# 欠采样：Tomek Link方法



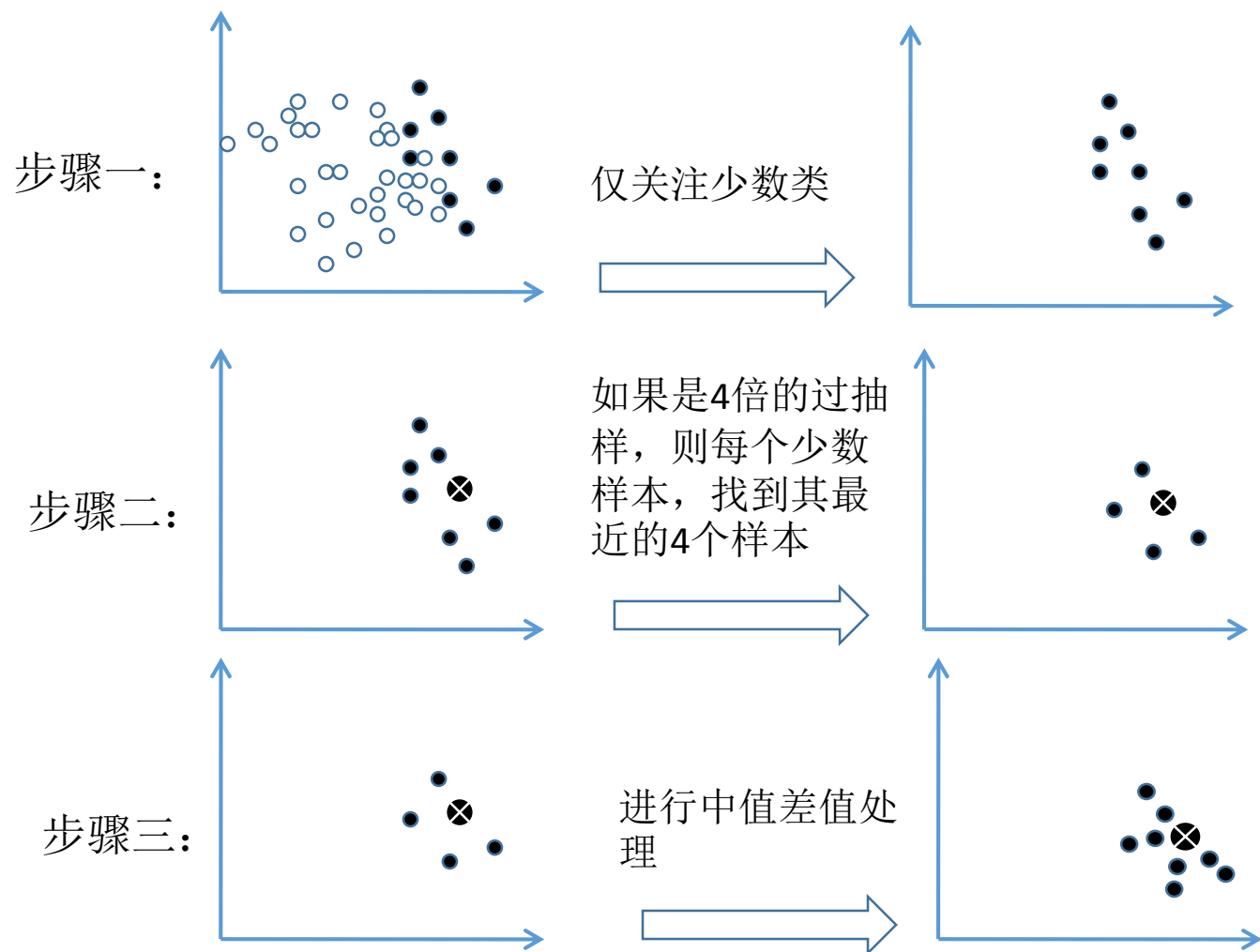
# 过采样：SMOTE方法的实现



# 过采样：SMOTE方法的问题

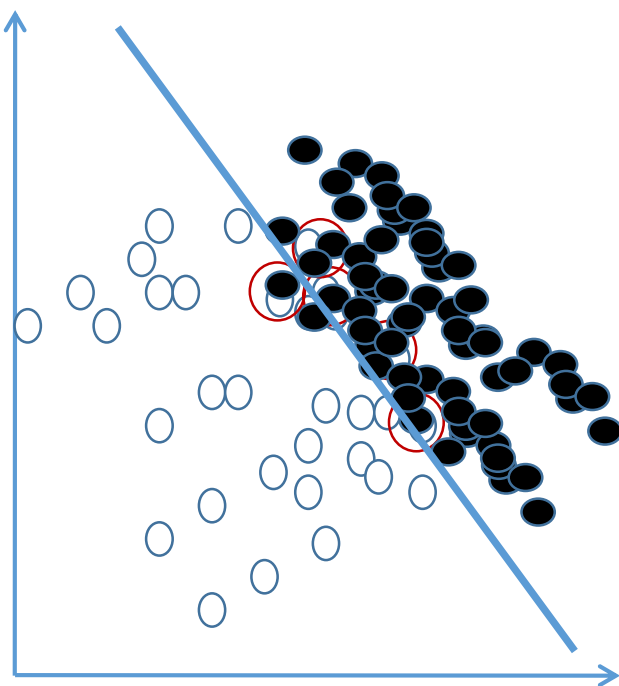


# 综合采样

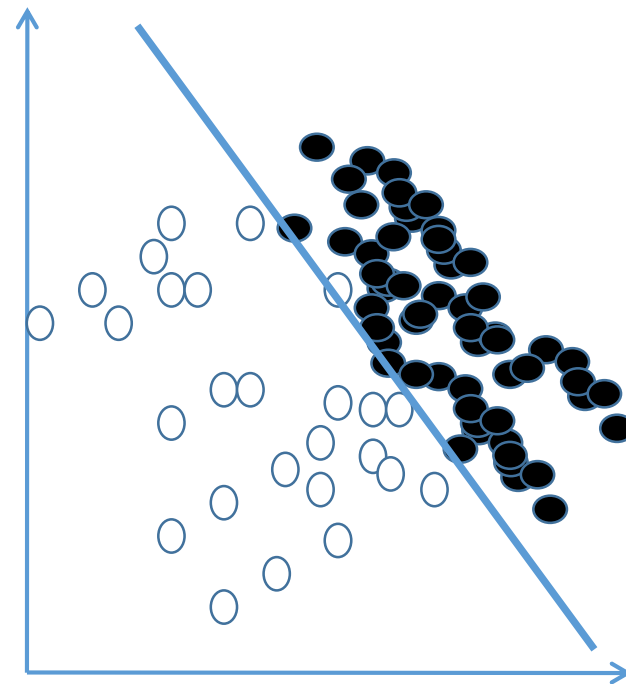


# 综合采样

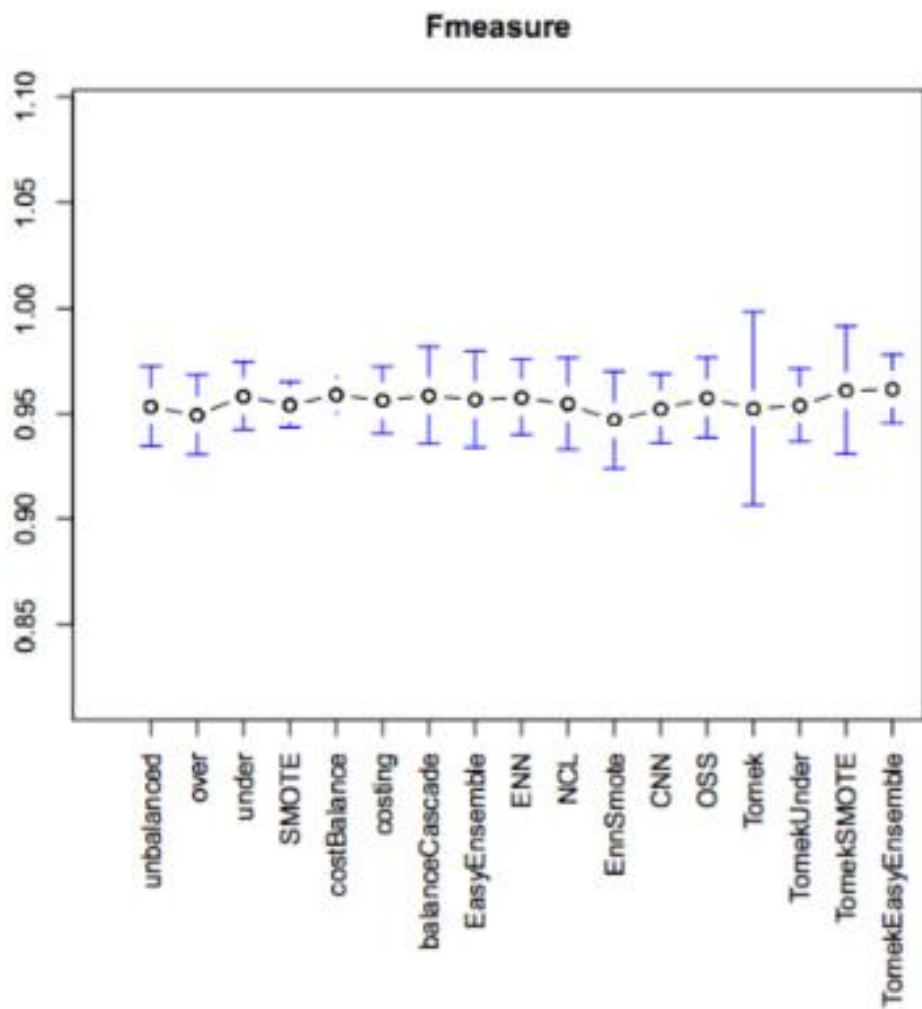
步骤四：



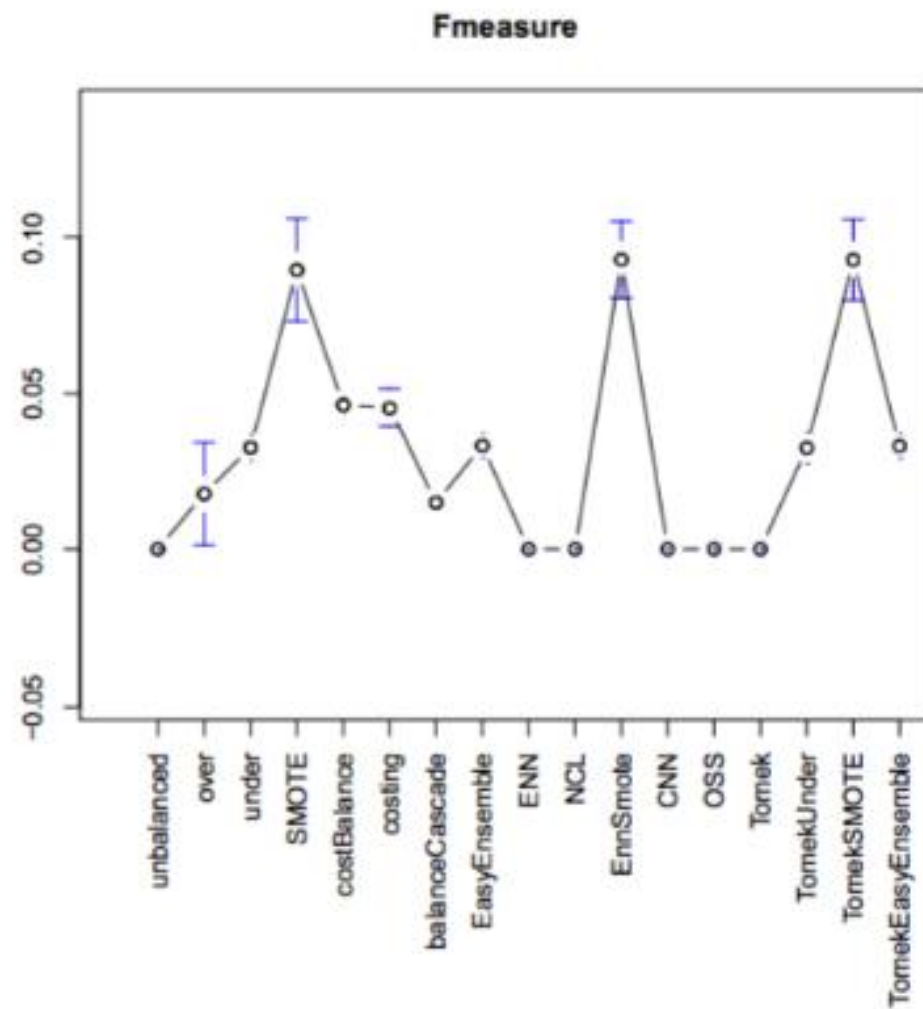
删除样本距离近的  
观测



# 只有检验才能知道



UCI breast cancer dataset



Atos fraud dataset

秦路主讲

## 七周成为数据分析师

七周为期，Get一条数据分析师职业黄金通道！



Python

## 数据分析与挖掘

集Python爬虫、数据采集、数据处理、数据分析与数据挖掘于一体，打造Python全栈工程师

主讲老师: 韦玮

VIP会员群+在线答疑+录播复习+1年反复观看

## 案例为师,实战为王

### 开启Python机器学习之路

科学规划全套课程体系,从入门到进阶,从理论到技巧,嵌入丰富课程案例讲解,逐步推进

讲师: 唐宇迪 深度学习领域多年一线实践研究专家

## 独一无二的 数据仓库建模指南系列教程升级版

- 从企业视角进行数据规划以及数据仓库模型的搭建
- 高质量的数据库模型和技巧,以及丰富的例子
- 数据仓库架构理论和实践要领

资深讲师: BAO胖子 15年+BI从业经验  
涉足电力、快消品、医药、信息服务行业的BI老兵

## 业务知识一站通

技术+业务,挣钱有门路!

讲师: 陈文



自己动手 丰衣足食

## Python3网络爬虫实战案例

一循序渐进,案例为王,诠释全面,思路制胜一

讲师: 崔庆才 北航硕士,百万级热度爬文博主



讲师 丘祐玮

## 人人都爱数据科学家

Python数据科学精华实战课程

## 数据分析报告制作

秘籍升级版

讲师: 陈丹奕 知乎大神,前百度资深数据分析师

## 先机致胜 破冰AI

深度学习模型/框架与实战

讲师: 唐宇迪 同济大学硕士  
深度学习领域多年一线实践研究专家



BI、商业智能  
数据挖掘 大数据  
数据分析师  
R语言 Python  
机器学习  
深度学习  
人工智能  
Hive Hadoop  
Tableau  
BIEE ETL  
数据科学家  
PowerBI