

# 深入浅出统计学

## 信息图形化

统计：通过某种有意义的方式对原始事实和数字进行提炼，使得仅仅通过观察原始数据无法立即水落石出的一些理念得以昭示

统计量是样本的函数

在样本未知的情况下是随机变量

对统计的研究：

- 统计数据来源
- 计算方法
- 有效使用方法得出结论

数据是指所搜索的原始事实与数字

信息是指加上某种意义的数据

- 饼图
- 条形图
  - 堆积条形图
  - 分段条形图
- 直方图
  - 长方形的面积与频数成比例
  - 图上的长方形之间没有间隔
  - 分组数值型数据
- 累积频数图
- 折线图

## 集中趋势的度量

均值： $\mu = \frac{\sum x}{n}$

中位数

- 数据右倾：均值>中位数
- 数据左倾：均值<中位数

众数

## 分散性与变异性的度量

### 分散性

极差/全距 = 上界 - 下界

迷你距：忽略异常值

- 四分位数
  - 最小四分位数 下四分位数 第一四分位数
  - 最大四分位数 上四分位数 第三四分位数
  - 中间的四分位数 中位数
- 四分位距 = 上四分位数 - 下四分位数
  - 与全距相比较少受到异常值的影响
  - 得到一种对几个数据集进行比较且比较结果不会被异常值扭曲的方法
- 下四分位数的位置
  - 首先计算  $\frac{n}{4}$
  - 如果结果为整数，则下四分位数位于  $\frac{n}{4}$  这个位置和下一个位置的中间

- 如果不是整数，则向上取整
- 上四分位数的位置
  - 首先计算  $\frac{3n}{4}$
  - 如果结果为整数，则下四分位数位于  $\frac{3n}{4}$  这个位置和下一个位置的中间
  - 如果不是整数，则向上取整
- 百分位数---将数据一分为百的数值
  - 第  $k$  百分位数就是位于数据  $k\%$  处的数值， $P_k$
  - 求百分位数
    - 首先将数值按升序排序
    - 为了求出  $n$  个数字的第  $k$  百分位数的位置，先计算  $k \left( \frac{n}{100} \right)$
    - 如果结果为整数，则百分位数处于第  $k \left( \frac{n}{100} \right)$  位和下一位数之间。取这两个位置上的数字平均值，得出百分位数
    - 如果不是整数，则将其向上取整，结果即为百分位数的位置
- 箱线图
  - 全距
  - 四分位距
  - 中位数

### 变异性

方差:  $\sigma^2 = \frac{\sum (x - \mu)^2}{n} = \frac{\sum x^2}{n} - \mu^2$

标准差:  $\sigma = \sqrt{\text{方差}}$

### 概率计算

- 概率是度量某事发生几率的一种数量指标
- 事件: 有概率可言的一个结果或一件事
- $P(A) = \frac{n(A)}{n(S)}$ 
  - $S$  样本空间
  - 事件是  $S$  的子集
- 对立事件
  - $P(A') = 1 - P(A)$
- 互斥事件
  - $P(A \cap B) = 0$
- 相交事件
- $P(A \cup B) = P(A) + P(B) - P(A \cap B)$
- 条件概率
  - $P(A|B) = \frac{P(A \cap B)}{P(B)}$
  - $P(A \cap B) = P(A|B)P(B)$
  - $P(B \cap A) = P(B|A)P(A)$
- 全概率公式
  - $P(B) = P(A)P(B|A) + P(A')P(B|A')$
- 贝叶斯定理
  - $P(A|B) = \frac{P(A)P(B|A)}{P(A)P(B|A) + P(A')P(B|A')}$
- 相关事件
- 独立事件
  - $P(A|B) = P(A)$
  - $P(A \cap B) = P(A)P(B)$

### 离散概率分布的运用

- 期望:  $E(X) = \sum xP(X = x)$ 
  - $E(X) = \mu$
- 方差:  $Var(X) = E(X - \mu)^2 = E(X^2) - \mu^2$

- 标准差:  $\sigma = \sqrt{Var(X)}$
- $E(aX + b) = aE(X) + b$
- $Var(aX + b) = a^2 Var(X)$
- 独立观察值:
  - $E(X_1 + X_2 + \dots + X_n) = nE(X)$
  - $Var(X_1 + X_2 + \dots + X_n) = nVar(X)$
- $X$ 和 $Y$ 是相互独立的随机变量
  - $E(X + Y) = E(X) + E(Y)$
  - $Var(X + Y) = Var(X) + Var(Y)$
  - $E(X - Y) = E(X) - E(Y)$
  - $Var(X - Y) = Var(X) + Var(Y)$

## 排列与组合

- $n$ 个对象排列数目:  $n! = n * (n - 1) * \dots * 3 * 2 * 1$
- $n$ 个对象作圆形排位:  $(n - 1)!$
- 按类型排位:
  - 第一类对象 $k$ 个, 第二类对象 $j$ 个, 第三类对象 $m$ 个...
  - $\frac{n!}{j! k! m! \dots}$
- 排列: 与顺序有关
  - $P_n^r = \frac{n!}{(n - r)!}$
- 组合: 与顺序无关
  - $C_n^r = \frac{n!}{r! (n - r)!}$

## 几何分布、二项分布及泊松分布

- 几何分布
  - $X$ 表示取得第一次成功所需进行的试验次数
  - $P(X = x) = q^{x-1}p$
  - $P(X > x) = q^x$
  - $X \sim Geo(p)$
  - $E(X) = \frac{1}{p}$
  - $Var(X) = \frac{q}{p^2}$
- 二项分布
  - $X$ 表示 $n$ 次试验中的成功次数
  - $P(X = x) = C_n^x p^x q^{n-x}$
  - $X \sim B(n, p)$
  - $E(X) = np$
  - $Var(X) = npq$
- 泊松分布
  - 泊松分布满足以下条件
    - 单独事件在给定区间内随机、独立的发生
    - 已知该区间内的事件平均发生次数
  - $X \sim Po(\lambda)$
  - $E(X) = \lambda$
  - $Var(X) = \lambda$
  - $P(X = x) = \frac{e^{-\lambda} \lambda^x}{x!}$
- 泊松分布近似二项分布
  - $\lambda \sim np$
  - $\lambda \sim npq$
  - 当 $n$ 很大且 $p$ 很小时, 可以用 $X \sim Po(np)$ 代替 $X \sim B(n, p)$

## 正态分布的运用

## 正态分布

- $X \sim N(\mu, \sigma^2)$

## 再谈正态分布的运用

### 二项分布的近似

如果  $X \sim B(n, p)$ , 且  $np > 5$ ,  $nq > 5$ , 则可以使用  $X \sim N(np, npq)$  近似代替二项分布

### 泊松分布的近似

如果  $X \sim Po(\lambda)$  且  $\lambda > 15$ , 则可用  $X \sim N(\lambda, \lambda)$  进行近似

## 统计抽样的运用

- 总体
  - 指的是对其进行测量、研究或分析的整个群体
- 样本
  - 从样本中选取的一部分对象
- 抽样单位
- 抽样空间
  - 列出总体中的所有独立单位
- 抽样方法
  - 简单随机抽样
    - 重复抽样
    - 不重复抽样
  - 分层抽样
  - 整群抽样
  - 系统抽样

## 总体和样本的估计

- 点估计量
  - 可用于估计总体参数数值的某个函数或算式
- 样本均值
  - $\bar{x} = \frac{\sum x}{n}$
- 样本均值估计总体均值
  - $\hat{\mu} = \bar{x}$
- 估计总体方差
  - $\hat{\sigma}^2 = \frac{\sum (x - \bar{x})^2}{n - 1}$
- 如果从一个非正态总体  $X$  中取出一个样本, 且样本很大, 则  $\bar{X}$  的分布近似为正态分布
  - $\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$

## 置信区间的构建

构建置信区间的步骤:

- 选择统计量
- 求出其抽样分布
- 决定置信水平
- 求出置信上下限

## 假设检验的运用

- 确定假设
  - 原假设即你要对其进行检验的断言, 除非有足够的证据进行反驳, 否则你将接受这个断言( $H_0$ )

- 备选假设即在拒绝 $H_0$ 之后将接受的另一个断言( $H_1$ )
- 选择检验统计量
  - 根据原假设 $H_0$ 选择检验统计量
- 确定拒绝域
  - 拒绝域分界点称为临界值
  - 为求拒绝域，先定显著性水平
    - 显著性水平度量的是一种愿望，即希望在样本结果的不可能程度达到多大时，就拒绝原假设 $H_0$
    - 显著性水平通常用 $\alpha$ 表示， $\alpha$ 越小，为了拒绝 $H_0$ ，样本结果需要达到的不可能程度越高
  - 单尾检验
  - 双尾检验
- 求出 $p$ 值
  - $p$ 值即为取得样本中的各种结果或取得拒绝域方向上的某些更为极端的结果的概率
- 样本结果是否位于拒绝域内
- 作出决策

第一类错误：错误地拒绝真原假设

- $P(\text{第一类错误}) = \alpha$

第二类错误：错误地接受假原假设

- 检查是否拥有 $H_1$ 的特定数值
- 求检验拒绝域以外的数值范围
- 假定 $H_1$ 为真，求得到这些数值的概率
- $P(\text{第二类错误}) = \beta$

假设检验的功效：在为 $H_0$ 假的情况下拒绝 $H_0$ 的概率

- 功效 =  $1 - \beta$