

多元统计分析—R 与 Python 的实现

回归

吴喜之

June 28, 2019

经典线性回归的数学假定:

- ① 模型对于参数的线性形式为:

$$y = X\beta + \epsilon \text{ 或 } y_i = \beta_0 + \beta_1 X_{i1} + \cdots + \beta_p X_{ip} + \epsilon_i, i = 1, 2, \dots, n;$$

- ② y 的所有值都是互相独立的, 如果把 X 看成是固定的, 则等价于 ϵ 的所有值都是互相独立的, 这意味着没有自相关;
- ③ 对于每个 X 的值, y 的分布是正态分布的, 如果把 X 看成是固定的, 则等价于 ϵ 的分布是正态分布;
- ④ 条件期望 $E(y|X) = X\beta$ 或 $E(\epsilon|X) = 0$, 协方差阵 $\text{Cov}(y|X) = \text{Cov}(\epsilon|X) = \sigma^2 I$, 这意味着所有观测值独立同方差 (homoscedasticity), 对于正态分布意味着所有观测值都独立同分布 (当然误差项也独立同分布);
- ⑤ 自变量和 ϵ 独立;

注意: 如果假定自变量是弱外生的 (**weak exogeneity**), 可以近似地把 X 看成是固定的, 上面的概率就不是条件概率了, 自变量和 ϵ 独立的条件就用不着了.

当自变量被看成是固定的, 人们所说的最小二乘线性回归的基本假定就是一句话: 观测值独立同正态分布. 或者用公式表示为: $y \sim N(X\beta, \sigma^2 I)$, $\epsilon \sim N(0, \sigma^2 I)$.

- ① 所有上述假定都无法核对.
- ② 把因变量和自变量的关系假定为线性关系的主要原因是当年人们的能力只能勉强应对线性模型所带来的数学及计算等一些问题.
- ③ 把样本假定为独立同正态分布, 这大多不合乎实际, 但基于与上一款同样的理由, 这种假定在数学上是比较方便的. 这种分布假定和最小二乘回归方法本身无关, 仅仅和与回归相关的各种推断有关, 比如 t 检验及 F 检验统计量的分布、各种估计量的分布及性质等等.

没有假定则无法做任何经典统计推断 (诸如假设检验等), 只有基于计算机算法的机器学习出现之后, 才有可能摆脱主观假定的束缚.