

回归中的相关度和决定系数

Pearson Correlation Coefficient.

(皮尔逊相关系数)

用于衡量两个值的线性相关的强度. $r \in (-1, 1)$.

r : 正向相关 > 0 , 负向相关 < 0 , 无相关性 $= 0$.

$$\rho = \text{Cor}(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X) \text{Var}(Y)}}$$

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

R^2 (决定系数)

若控制自变量不变, 则因变量的变异程度会减少 R^2 .

如: $R^2 = 0.8$

即若自变量 $x = x_0$ 则有 80% 的 $y = \hat{y}$ 而 20% 的 $y \neq \hat{y}$.

对于 simple linear regression

$$R^2 = r \cdot r.$$

对于 multiply regression.

对于 multiply regression.

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

注: SSR : $\sum_{i=1}^n (\hat{y}_i - \bar{y})^2$
↓ regression
sum sqrt 回归模型
平方和和 由于模型的建立而
引起的差异.

SST : $\sum_{i=1}^n (y_i - \bar{y})^2$
↓ total
总和. 点的总变异量.

SSE : $\sum_{i=1}^n (y_i - \hat{y}_i)^2$
↓ error
误差值 点的误差值.

$$SSR + SSE = SST.$$

局限性: R^2 会随着自变量的增多而变大.

修正:

$$R^2_{\text{adjusted}} = 1 - \frac{(1 - R^2)(N - 1)}{N - P - 1}$$

↓ 修正. ↗ 样本误差.
↘ 预测量误差.