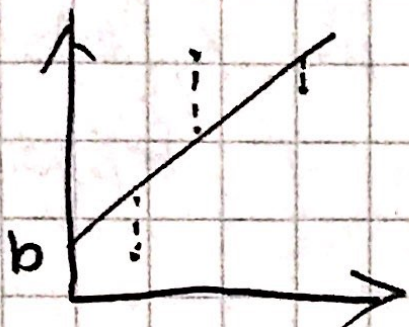


线性回归中的平方误差及公式推导



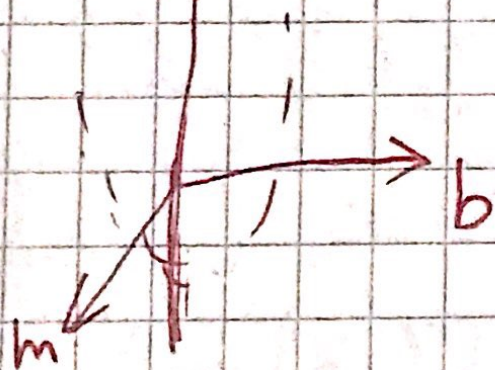
$$SE_{line} = (y_1 - (mx_1 + b))^2 + \dots + (y_n - (mx_n + b))^2$$

$$= \sum y_n^2 - 2m \sum x_n y_n - 2b \sum y_n$$

$$= \sum y_n^2 + m^2 \sum x_n^2 + 2mb \sum x_n + nb^2$$

$$= n \bar{y}^2 + 2m \bar{x} \bar{y} - 2b n \bar{y} + m^2 n \bar{x}^2 + 2mnb \bar{x} + nb^2$$

\bar{x}, \bar{y} 可视为已知, 只有 m, b 未知



求导
↓

$$m = \frac{\overline{xy} - \bar{x}\bar{y}}{(\overline{x^2} - \bar{x}^2)}$$

决定系数 R^2

y 的波动多大程度可被 x 描述的波动

$$R^2 = 1 - \frac{\text{SE}_{\text{line}}}{\text{SE}_{\bar{y}}}$$

离 总波动

$$\text{SE}_{\bar{y}} = \sum (y_n - \bar{y})^2$$

$$\text{SE}_{\text{line}} = \sum (y_n - (mx + b))^2$$

协方差 & 回归线

covariance

$$\text{Cov}(X, Y) = E[(X - E[X])(Y - E[Y])]$$

与最小二乘联系 $m = \frac{\text{Cov}(X, Y)}{\text{Var}(X)}$

$$E[(X - E[X])(Y - E[Y])]$$

$$= E[XY - XE[Y] - E[X]Y + E[X]E[Y]]$$

把 $E[X]$ 、 $E[Y]$ 看作常数

$$E[XY] - E[X]E[Y] - E[X]E[Y] + E[X]E[Y]$$

$$= \underbrace{E[XY]}_{\approx \overline{XY}} - \underbrace{E[X]E[Y]}_{\overline{Y}\overline{X}}$$

$$\text{Cov}(X, Y) = E[XY] - E[X]E[Y]$$

$(= \overline{XY} - \overline{X}\overline{Y})$ 回归的分子

χ^2 分布

$$X \sim N(0, 1)$$

$$Q = X^2$$

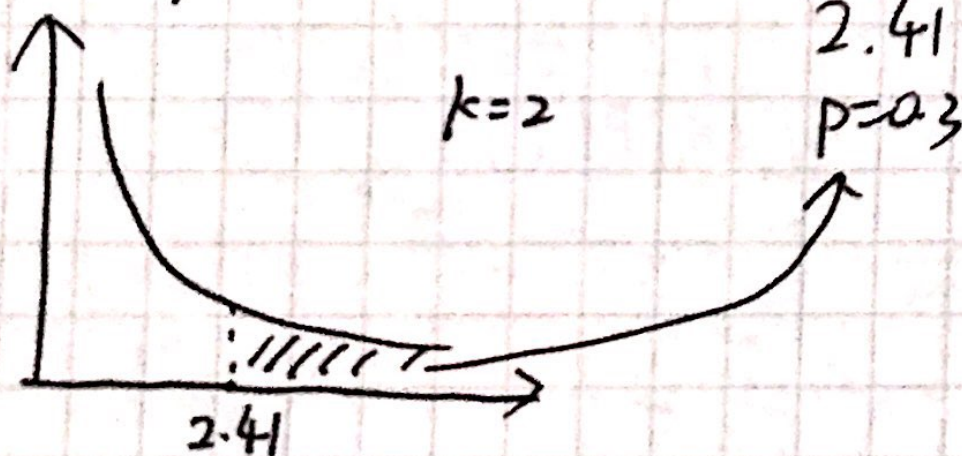
$$\rightarrow \cancel{X^2} \quad Q \sim \chi_{df}^2$$

χ^2 分布

$$X \sim N(0, 1)$$

$$Q = X_1^2 + \dots + X_k^2$$

$$\downarrow Q \sim \chi_k^2$$



皮尔逊 χ^2 检验

$$\chi^2 = \sum \frac{(f - f_e)^2}{f_e}$$
$$df = n - 1$$

列联表 χ^2 检验

Sick	A	B	C	80
Nt Sick	D	E	F	300
				380

(21%)
(79%)

$$\chi^2 = \sum \frac{(f - f_e)^2}{f_e}$$

$$df = (行 - 1)(列 - 1)$$