

Setup

$$y = X\beta + \varepsilon$$

given obs y we have

fitted values: $\hat{y} = X\hat{\beta}$

estimated residuals: $\hat{\varepsilon} = y - \hat{y} = y - X\hat{\beta}$

where $\hat{\beta} = (X^T X)^{-1} X^T y$

DEF

The hat matrix is

$$H = X(X^T X)^{-1} X^T$$

note this idempotent:

$$H^2 = H \cdot H = X(X^T X)^{-1} X^T X (X^T X)^{-1} X^T = X(X^T X)^{-1} X^T = H$$

The fitted values are

$$\hat{y} = X \hat{\beta} = X(X^T X)^{-1} X^T y = Hy$$

So H projects data onto fitted values.

Estimated resid are

$$\hat{\varepsilon} = y - \hat{y} = y - Hy = (I - H)y$$

Thus, $\boxed{\text{Var } \hat{\varepsilon} = \sigma^2 (I - H)}$ [proof: HW]

So, the SE for $\hat{\varepsilon}_i$ is

$$SE(\hat{\varepsilon}_i) = \hat{\sigma} \sqrt{1 - H_{ii}}$$

 ith element of $\text{diag}(H)$

Under A1, the studentized residuals are

$$\hat{\varepsilon}_i^* = \frac{\hat{\varepsilon}_i - 0}{SE(\hat{\varepsilon}_i)} \sim N(0, 1)$$

which gives us a quantitative way to assess outliers, i.e., if $\hat{\varepsilon}_k^*$ is outside of $[-2, 2]$ then it is potentially an outlier.

[in R "standardized residual" = "studentized residual"]

Note

$$\hat{y} = Hy$$

$$\Rightarrow \hat{y}_i = h_{i1}y_1 + \dots + h_{ii}y_i + \dots + h_{in}y_n$$

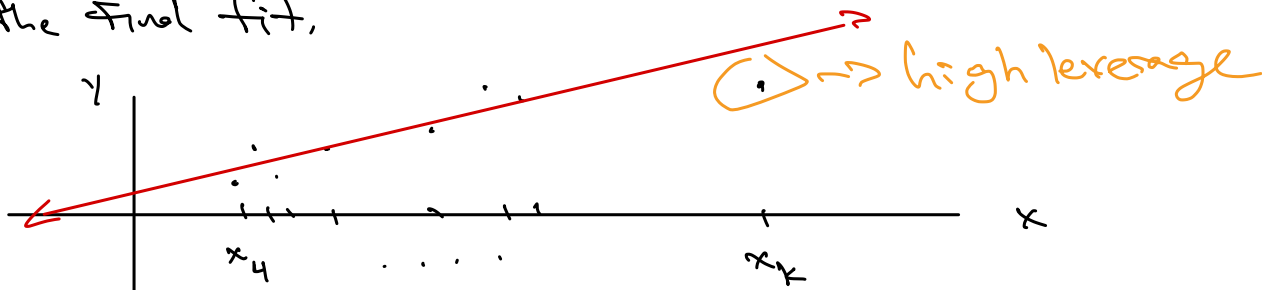
so h_{ii} is the influence of y_i on its own fitted value \hat{y}_i

DEF

Define h_{ii} to be the leverage of (y_i, x_i)

Note

Leverage is the potential for y_i having an influence on the final fit,



data points with unusual x values have high leverage

Remark $H = X(X^T X)^{-1} X^T$ only depends on X , not the actual value of y_i , so it only measures potential influence.

To measure actual influence, Cook's Distance (Cook's D) is

$$D_i = \frac{H_{ii}}{p(1-H_{ii})} (\hat{\varepsilon}_i^*)^2$$

leverages are less than 1

Note

High leverage \Rightarrow potential for influencing fit

High Cook's $D \Rightarrow$ actual influence over fit

[R]

Assessing quality of model fit

How do we compare competing models, i.e., assess goodness-of-fit?

Recall

$$RSS = (y - X\hat{\beta})^T (y - X\hat{\beta}) = \sum_i (y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \dots + \hat{\beta}_p x_{ip}))^2$$

and

$$r^2 = 1 - \frac{RSS}{SYY} \quad \text{where } SYY = \sum (y_i - \bar{y})^2$$

Note $r^2 \in [0, 1]$, larger values \Rightarrow better fit

RSS, smaller values \Rightarrow better fit.

But: Both will always improve with the introduction of more features.

If model has d features, define

Mallow's C_p : $C_p = \frac{1}{n} (RSS + 2d\hat{\sigma}^2)$

Akaike information criterion (AIC) : $AIC = \frac{1}{n\hat{\sigma}^2} (RSS + 2d\hat{\sigma}^2)$

Bayesian information criterion (BIC) : $BIC = \frac{1}{n} (RSS + \log(n)d\hat{\sigma}^2)$

We prefer models with smaller C_p / AIC / BIC.

BIC tends to choose simpler models than AIC.