

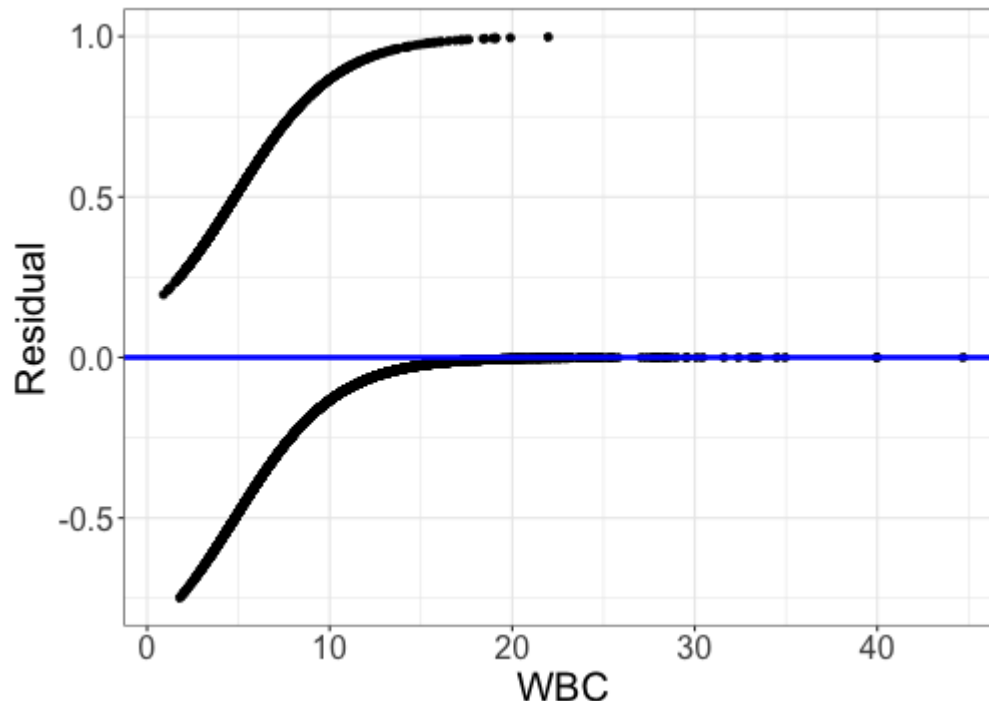
Logistic regression assumptions and diagnostics

(residual plots $y - \hat{p}$)

Don't use usual residuals for logistic regression

Fitted model: $\log\left(\frac{\hat{p}_i}{1 - \hat{p}_i}\right) = 1.737 - 0.361 WBC_i$

Residuals $Y_i - \hat{p}_i$:



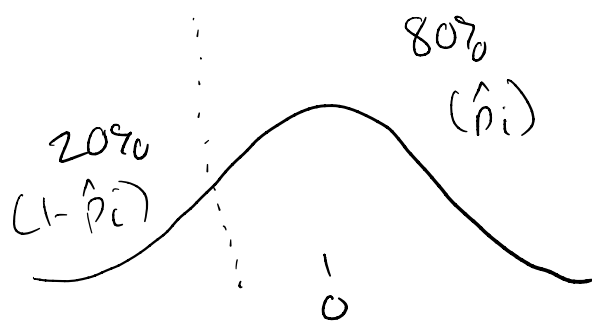
(checking shape assumption)
(randomized)

Quantile residuals for logistic regression

Motivation : Suppose $\hat{p}_i = 0.8$. I want to create residual r_Q that behaves like linear regression residuals : want

- If $\hat{p}_i \approx p_i$ (good estimate) then $E[r_Q | x_i] \approx 0$
- If $\hat{p}_i > p_i$ (overestimate), then $E[r_Q | x_i] < 0$
- If $\hat{p}_i < p_i$ (underestimate), then $E[r_Q | x_i] > 0$
- want $r_Q \approx \text{Normal}$ (if $\hat{p}_i \approx p_i$)

Idea : $\hat{p}_i = 0.8$, Divide $N(0, 1)$ into 2 regions:



- If $y_i = 1$, sample r_Q from right side
- If $y_i = 0$, sample r_Q from left side

If $\hat{p}_i \approx p_i$, then I'm sampling from a $N(0, 1)$ on average

Pseudo-code:

for each $i = 1, \dots, n$:

Calculate \hat{p}_i

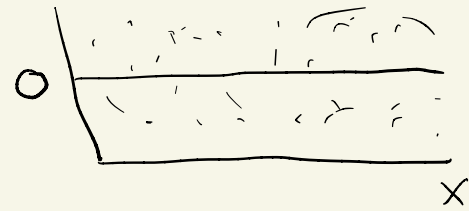
if $Y_i = 1$: once

Sample Y from the upper \hat{p}_i area of $N(0,1)$

if $Y_i = 0$

Sample once from the lower $1 - \hat{p}_i$ area of $N(0,1)$

use Q to make residual plots
quantile
residuals



If $\hat{p}_i \approx p_i$, then $Q_i \sim N(0,1)$ (over many datasets)

If all $\hat{p}_i \approx p_i \forall i$, then marginally $Q \sim N(0,1)$

Class activity, Part I

https://sta711-s23.github.io/class_activities/ca_lecture_10.html

Leverage and Cook's distance

Linear regression : $\hat{y} = X\hat{\beta}$
 $= X \underbrace{(X^T X)^{-1} X^T}_{\text{"hat matrix" } H} y$

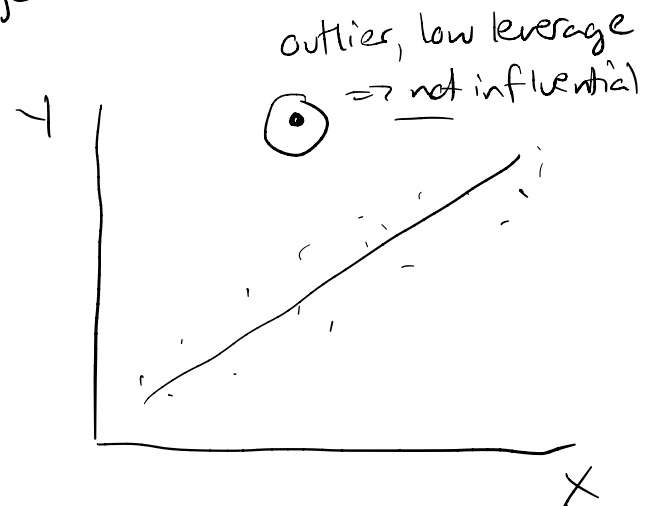
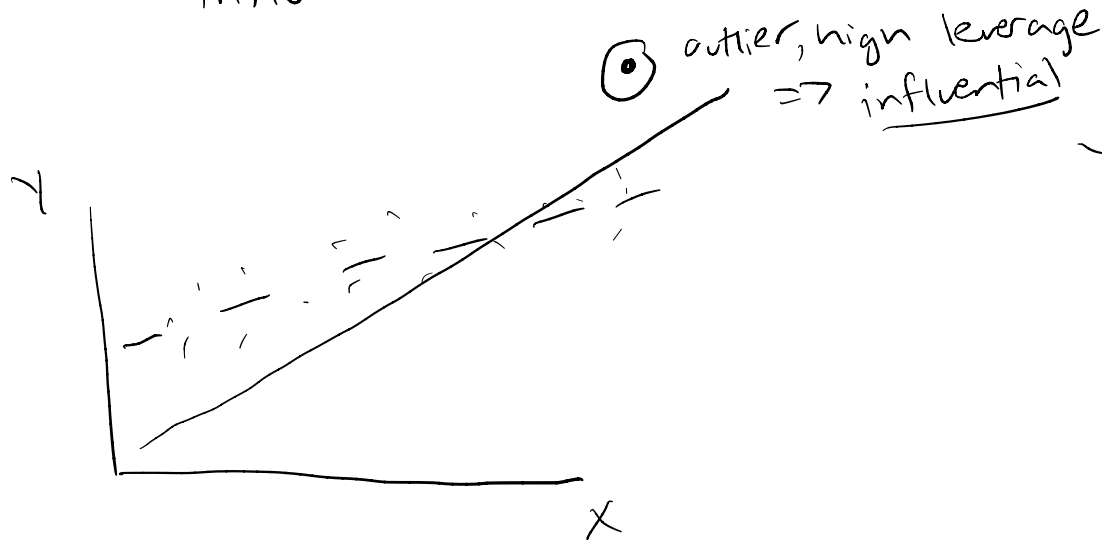
(not Hessian! I will
 $H(\beta)$ to denote Hessian)

$$\text{Var}(y - \hat{y}) = \sigma^2(I - H)$$

$$\Rightarrow \text{Var}(y_i - \hat{y}_i) = \sigma^2(1 - h_i)$$

leverage = potential of observation to influence fit

$h_i = [H]_{ii}$ leverage of the i^{th} observation



Cook's distance (linear regression) :

$$D_i = \frac{(y_i - \hat{y}_i)^2}{\underbrace{(n+1) \hat{\sigma}^2}_{\text{\# of } \beta \text{ in model}}} \cdot \frac{h_i}{\underbrace{(1-h_i)^2}_{\text{high leverage?}}}$$

outlier?

concerned that a point is influential when $D_i > \text{threshold}$ (e.g. 0.5 or 1)

Logistic regression :

$$(W = \text{diag}(p_i(1-p_i)))$$

$$\text{Hat matrix } H = W^{\frac{1}{2}} X (X^T W X)^{-1} X^T W^{\frac{1}{2}}$$

$h_i = \text{leverage}$

$$D_i = \frac{(y_i - \hat{p}_i)^2}{(n+1) \hat{p}_i (1 - \hat{p}_i)} \cdot \frac{h_i}{(1-h_i)^2}$$

concerned when $D_i > 0.5$ or 1

Class activity, Part II

https://sta711-s23.github.io/class_activities/ca_lecture_10.html