# Outlier Detection

**Model:**
$$Y_i \mid \beta, \sigma^2 \overset{ind}{\sim} N(x_i^T \beta, \sigma^2)$$

**Residual Analysis:**
$$\frac{Y_i - x_i^T \beta}{\sigma} \mid \beta, \sigma^2 \overset{iid}{\sim} N(0, 1)$$

Need to estimate
$$\hat{e}_i = \frac{Y_i - x_i^T \hat{\beta}}{\hat{\sigma}} \overset{\cdot}{\sim} N(0, 1)$$

$\uparrow$
if good
estimates

Should be rare to see $|\hat{e}_i| \geq 2$ so detect outliers with the rule

$$\left| \frac{Y_i - x_i^T \hat{\beta}}{\hat{\sigma}} \right| > K_\alpha \quad \text{for} \quad K_\alpha = -\Phi(\alpha/2)$$

In the Bayesian paradigm $\beta \mid Y, \sigma^2 \sim N(\hat{\beta}, \Sigma_\beta)$

So $Y_i - x_i^T \beta \mid Y, \sigma^2 \sim N(Y_i - x_i^T \hat{\beta}, \ x_i^T \Sigma_\beta x_i)$

Let's assume for simplicity $P(\beta, \sigma^2) \propto \frac{1}{\sigma^2}$

$$\beta \mid \sigma^2, Y \sim N(\hat{\beta}_{MLE}, \ \sigma^2 (X^T X)^{-1})$$

$$\sigma^2 \mid Y \sim \text{Inv-}\chi^2 \left( n-p, \ \frac{RSS(\hat{\beta}_{MLE})}{n-p} \right) \equiv \text{Inv-}\chi^2 \left( n-p, \ \hat{\sigma}^2_{MLE} \right)$$

$\beta | y \sim t_{n-p} \left( \hat{\beta}_{MLE}, \hat{\sigma}^2 (X^T X)^{-1} \right)$

With this, it makes sense to discuss

$$Y_i - x_i^T \beta \,|\, \sigma^2, y \sim N\left( Y_i - x_i^T \hat{\beta}_{MLE}, \sigma^2 \underbrace{x_i^T (X^T X)^{-1} x_i}_{H_{ii}} \right)$$

$$Y_i - x_i^T \beta \,|\, y \sim t_{n-p} \left( y_i - x_i^T \hat{\beta}, \; \hat{\sigma}^2 h_{ii} \right)$$

For large $n$, $\sigma^2 | y$ concentrates around $\hat{\sigma}^2 = \dfrac{RSS(\hat{\beta}_{MLE})}{n-p}$ and
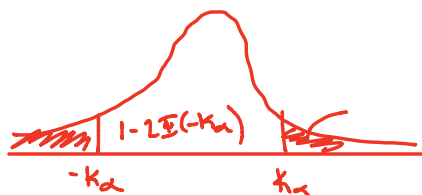
$$\frac{Y_i - x_i^T \beta}{\sigma} \,|\, y \,\dot{\sim}\, N\left( \frac{y_i - x_i^T \hat{\beta}}{\hat{\sigma}}, \; h_{ii} \right)$$

Posterior Residual dist.

With this we can calculate $\mathbb{P}\left( \left| \frac{y_i - x_i^T \beta}{\sigma} \right| \leq k_\alpha \,\Big|\, Y \right)$ will do this numerically.

We will need a multiple test correction.

Let $p(k_\alpha) = $ prob of being an outlier then we want

$$\mathbb{P}(\text{no outliers}) = p(k_\alpha)^n \equiv 1 - \alpha$$



$1 - 2\Phi(-k_\alpha)$

$-k_\alpha \qquad k_\alpha$

$$p(k_\alpha) = 1 - 2\,\Phi(-k_\alpha) = (1-\alpha)^{1/n}$$

$$\implies \boxed{K_\alpha = - \mathcal{I}^{-1} \left( \frac{1}{2} \left( 1 - (1-\alpha)^{1/2} \right) \right)}$$  <u>Rmk:</u> Bonferroni Correction.

Finaly, we can use this for model checking.

$$Y - X\beta \mid \sigma^2, Y \sim N\left( Y - X\hat{\beta}, \sigma^2 \underbrace{X(X^TX)^{-1}X^T}_{H} \right)$$

$H$ may be rank deficent; $X = QR$, $\hat{\beta} = (X^TX)^{-1}X^TY$

$$= (R^TQ^TGR)^{-1} R^TQ^TY$$

$$= R^{-1}R^{-T}R^T Q^TY$$

$$\implies \boxed{R\hat{\beta} = Q^TY}$$  numerically stable to find $\hat{\beta}$

$$Q^T(Y - X\beta) \mid \sigma^2, Y \sim N\left( Q^TY - Q^TQR\hat{\beta}, \sigma^2 Q^TQ Q^TQ \right)$$

$$\equiv N\left( 0, \sigma^2 I_p \right)$$

$$\implies \frac{1}{\sigma^2} Q^T(Y - X\beta) \mid Y \sim N(0, I_p)$$

Suggests use of a QQ-plot on standardized residuals