

Week 3-02424

! Linearity assumption of the input! in GLM besides Gaussianity

$$\bullet N(\beta X, \sigma^2 I)$$

parameters
values

↑ Independence assumption
Design matrix

- col 1 = col 2 + col 3 \Rightarrow linear dependence, i.e. not structural identifiable
col 1 = col 4 + col 5 + col 6

(Slide 7)

(Same as singular X)

- fix:

	B_1	B_2	B_3
A_1	μ_0	$\mu_0 + \beta_2$	$\mu_0 + \beta_3$

A_2	$\mu_0 + \alpha_2$	$\mu_0 + \alpha_2 + \beta_2$	$\mu_0 + \alpha_2 + \beta_3$
-------	--------------------	------------------------------	------------------------------

- Distance: $\delta_Z(y_1, y_2) = y_1^T \Sigma^{-1} y_2 \Rightarrow \|y\|_Z = \sqrt{\delta_Z(y, y)}$
- Deviance: $D(y; \mu) = \delta_Z(y - \mu, y - \mu)$

- factor: discrete
- covariate: continuous

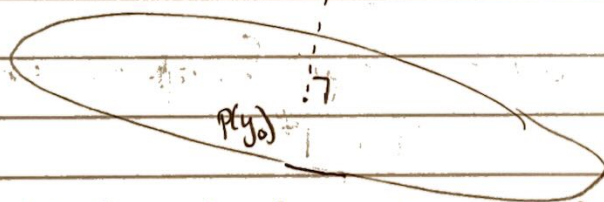
- Correlation in data which is not taken in consideration will lead to a too confident model.

- Very correlated data do not contain n information but pn information. Hence too much information is assumed

- orthogonal: $\sum (p_0(y); y - p_0(y)) = 0$ (slide 27)

- CI: $\hat{\beta}_i \pm t_{1-\alpha/2}^{n-p} \text{SE}_i$ fordi vi estimerer σ^2 i t fra slide 28. Vi tager højde for usikkerheden i estimatet, for σ^2 ved at bruge en t -fordeling i stedet for en N .

- Residuals are orthogonal to the model (slide 35)



Vi kan kun variere i de resterende dimensioner $n-p$ fordi de p dimensioner kan beskrives ved modellen.

- exercise slide 40

$$\left. \begin{aligned} y_1 &= \beta_A + \varepsilon_1 \\ y_2 &= \beta_A + \underbrace{(\beta_B - \beta_A)}_{\gamma} + \varepsilon_2 \\ y_3 &= 2\beta_A + \gamma + \varepsilon_3 \end{aligned} \right\} \Rightarrow X = \begin{bmatrix} 1 & 0 \\ 1 & 1 \\ 2 & 1 \end{bmatrix} \Rightarrow \text{Same } H \text{ as other design from slide 39}$$

under constraint at $y_i \sim N(0, 1) \quad \forall i$

• $Y^T Y = Y^T (I - H) Y + Y^T H Y$

$X = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$ $H = \begin{bmatrix} 1 & \dots & 1 \\ \vdots & & \vdots \\ 1 & \dots & 1 \end{bmatrix}$ under $X \sim H$ (Slide 44)

$\sum (y_i - \bar{y})^2 = \underbrace{\left(\sum y_i^2 \right)}_{\chi^2_n} - \underbrace{n \bar{y}^2}_{\chi^2_1}$

χ^2_{n-1}

• R^2_{adj} takes complexity into account because

R^2 for the full model has $R^2 = 1$ but is VERY complex

• Slide 50:

$$F(y) = \frac{D(p_1(y); p_0(y)) / (m_1 - m_0)}{D(y; p_1(y)) / (n - m_1)}$$

- Where $D(p_1(y); p_0(y))$ can be understood as $Y^T (H_0 - H_1) Y$ and distributed as $\chi^2_{m_1 - m_0}$

- And $D(y; p_1(y))$ can be understood as $Y^T (I - H_1) Y$ and distributed as $\chi^2_{n - m_1}$

- We then normalize with the degrees of freedom to obtain the F-distribution

i.e. we have the squared distances distributed by χ^2 -distribution and the more degrees of freedom the larger the variation.

In the F-test we then test if it is worth to go down in complexity. see slide 49.