

For podcast: technical problems with audio, hopefully solved.

Recap

- introduction
- linear model
- deviance

$$D(\hat{\beta}) = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \|y - \hat{y}\|^2$$

$$\hat{\sigma}^2 = s^2 = \frac{D(\hat{\beta})}{n-p}$$

$n = \# \text{ observations}$
 $p = \# \text{ variables}$

$$\text{Var}(\hat{\beta}) = s^2 (X^T X)^{-1}$$

under assumption of normality $\frac{\hat{\beta} - \beta_0}{\sqrt{s^2 (X^T X)^{-1}}} \stackrel{H_0}{\sim} \mathcal{N}(0; 1)$

\rightarrow t-value

$\Pr(\text{obtaining a t-value "worse" than that actually obtained}) = p\text{-value}$

small p-value \rightarrow evidence against H_0

- polynomial
- dummy variables for categorical variables I_A

Today

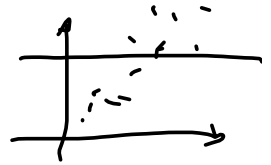
- R^2 & graphical diagnostic
- variable transformations
- multivariate response
- computational tricks.

To evaluate the goodness of fit, we need to compute the coefficient of determination

$$R^2 = \frac{\text{explained deviance}}{\text{total deviance}} = 1 - \frac{\text{residual deviance}}{\text{total deviance}}$$

$$= 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

$\bar{y} = \sum_{i=1}^n y_i$



R^2 : - simple to interpret
(fraction of variability explained by the model)

- oversimplification
(everything is reduced to a number) → graphical tools

⇒ graphical diagnosis, are based on the residuals

$$\hat{\epsilon}_i = y_i - \hat{y}_i \quad i = 1, \dots, n$$

$$\text{Var}(\epsilon_i) = \sigma^2$$

($\hat{\epsilon}_i$ are surrogate of ϵ_i , that are not observable)

- Figure 2.4(a) → Anscombe plot
check violation of homoscedasticity
- Figure 2.4(b) → quantile-quantile plot (qqplot) $\epsilon_i \sim N(0, \sigma^2)$
check violation of normality
y-axis = (standardized) values of $\hat{\epsilon}_i$
x-axis = expected value under normality

General conclusions

- model is decent, especially for the 'average' cars;
- simplicity;
- graphical diagnoses not totally satisfying
- model is bad for extrapolation;
- 'not realistic' (cannot increase) the distance → for larger engine sizes

Variable transformations

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \varepsilon$$

linear model

↳ refers to the parameters (linear in the parameters)

We can use whatever transformation of the variables, as long as the parameters have a linear relationship

$$y = \beta_0 + \beta_1 x_1 + \beta_2 \frac{1}{x_2} + \beta_3 \frac{x_1}{x_2} + \beta_4 e^{x_2^2 + x_1} + \varepsilon$$

We can also transform y

E.g. transformation of the response

$$\text{consumption} = \beta_0 + \beta_1 (\text{engine size}) + \beta_2 I_A + \varepsilon$$

↑
distance covered

+ nice all
+ $R^2 = 0.64$
+ polynomial terms are not necessary of higher orders
↳ simplify the interpretation

going back to original scale:

- not totally satisfactory (left ~~part of the plot~~ part of the plot including gas cars)
- $R^2 = 0.56$
- graphical diagnostics are unsatisfactory

Way more used transformation: logarithm

- especially good for variables with support $(0; +\infty)$, to transform into a variable with support \mathbb{R} (more in line with normality assumptions)

$$\log(\text{distance covered}) = \beta_0 + \beta_1 \log(\text{engine size}) + \beta_2 I_A + \varepsilon$$

We cannot do much better with these variables, but we have a lot of additional information to try to explain the variability of y

→ add new variables to the model!

- Weight, we know that can affect the distance per liter
- there are always points far from the others in the bottom left of the plots
→ they belong to cars with engines with only 2 cylinders
→ new dummy variables 2 cylinders / rest

$$I_D = \begin{cases} 1 & \text{if the engine has 2 cylinders} \\ 0 & \text{otherwise} \end{cases}$$

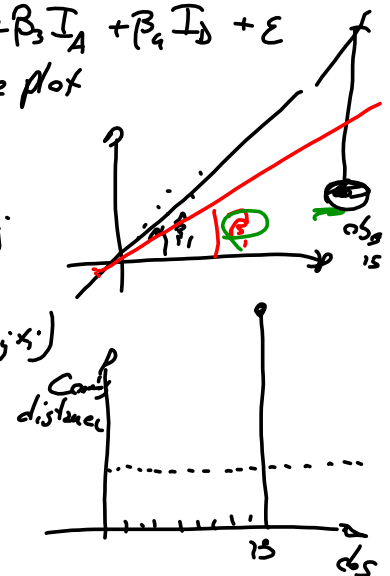
Model:

$$\log(\text{distance}) = \beta_0 + \beta_1 \log(\text{engine size}) + \beta_2 \log(\text{weight}) + \beta_3 I_A + \beta_4 I_D + \varepsilon$$

quite good result in Anscombe plot and quantile-quantile plot

Two additional graphical diagnostic tools:

- scatter plot of the residuals $\text{Cor}(\varepsilon_i, \varepsilon_j) = 0 \quad \forall i \neq j$
- Cook's distance
- effect on $\hat{\beta}$ when removing a specific observation (y_i, x_i)
- check if any point influence too much the estimates



Multivariate response

We have a multivariable model → more than one explanatory variable
multivariate model → more than one response variable

$$Y = (y_1, y_2, \dots, y_q)$$

$$Y = XB + E$$

$(n \times q) \quad (n \times p) (p \times q) \quad (n \times q)$

$$B = \begin{pmatrix} \beta_{11} & \beta_{12} & \dots & \beta_{1p} \\ \vdots & \vdots & \ddots & \vdots \\ \beta_{q1} & \beta_{q2} & \dots & \beta_{qp} \end{pmatrix}$$

$$\text{Var}(E) = \Sigma$$

$(q \times q)$

$\sum_{i=1}^q$
 $\Sigma = \text{covariance matrix}$

$$\hat{B} = (X^T X)^{-1} X^T Y$$

$$\hat{\Sigma} = \frac{1}{n-p} Y^T P Y$$

Computational aspects

→ very important in data analysis due to the size of the data (large n , large p)

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

↳ for the least squares estimator the most problematic task is the inversion of $X^T X$

→ Choleski factorization

A positive definite matrix

$$A = L L^{*T} \quad \text{where } L \text{ is a lower triangular matrix}$$

↑
unique

L^{*T} is the ~~conjugate~~ transposed

in \mathbb{R}

$$\begin{pmatrix} 4 & 12 & -16 \\ 12 & 37 & -43 \\ -16 & -43 & 98 \end{pmatrix} = \begin{pmatrix} 2 & 0 & 0 \\ 6 & 1 & 0 \\ -8 & 5 & 3 \end{pmatrix} \begin{pmatrix} 2 & 6 & -8 \\ 0 & 1 & 5 \\ 0 & 0 & 3 \end{pmatrix}$$

A

L

L^T

computational time $O(p^3 + np^2/2)$

→ Use the Gram-Schmidt process to compute the QR decomposition
orthogonalization

$$X = QR \quad \text{such that } Q \text{ } n \times p \text{ matrix, often } Q^T Q = I$$

R $p \times p$ " , upper triangular

then $\hat{\beta} = R^{-1} Q^T y$ → easy to invert because R is upper triangular

$$\hat{y} = Q Q^T y$$

$$\begin{aligned} \hat{\beta} &= (X^T X)^{-1} X^T y \\ &= (R^T Q^T Q R)^{-1} Q^T R^T y \\ &= (R^T)^{-1} \underbrace{I}_{\text{to solve}} R^{-1} Q^T y \end{aligned}$$

n very large \rightarrow storage problems
row by row procedure

$$\hat{\beta} = \underbrace{(X^T X)^{-1}}_{W = X^T X} \underbrace{X^T y}_U = X^T y$$

W are the only terms needed
 U to compute the OLS

$$X = \begin{pmatrix} \tilde{x}_1^T \\ \tilde{x}_2^T \\ \vdots \\ \tilde{x}_n^T \end{pmatrix} \quad \text{where } \tilde{x}_i^T \text{ is the } i\text{-th row of } X$$

$$\text{Then } W = \sum_{i=1}^n \tilde{x}_i \tilde{x}_i^T, \quad U = \sum_{i=1}^n \tilde{x}_i y_i$$

$$W_{(j)} = W_{(j-1)} + \tilde{x}_j \tilde{x}_j^T, \quad U_{(j)} = U_{(j-1)} + \tilde{x}_j y_j \quad j = 2, \dots, n$$

$$W_{(1)} = \tilde{x}_1 \tilde{x}_1^T, \quad U_{(1)} = \tilde{x}_1 y_1$$

Recursive estimation

- we have W , we still need to invert it $\frac{W}{p \times p} \leftarrow$ problematic to invert, especially for large p

$$\hat{\beta}_{(n)} = V_{(n)} X_{(n)}^T y$$

$$\hookrightarrow W_{(n)}^{-1} = (X_{(n)}^T X_{(n)})^{-1}$$

update with a new observation $n+1$ $(\tilde{x}_{n+1}, y_{n+1})$

$$X_{(n+1)} = \begin{pmatrix} X_{(n)} \\ \tilde{x}_{n+1}^T \end{pmatrix} \Rightarrow W_{(n+1)} = W_{(n)} + \tilde{x}_{n+1} \tilde{x}_{n+1}^T$$

$$= X_{(n)}^T X_{(n)} + \tilde{x}_{n+1} \tilde{x}_{n+1}^T$$

SHERMAN-MORRISON FORMULA (see appendix A.1, formula (A.2))

$$(A + b d^T)^{-1} = A^{-1} - \frac{1}{1 + d^T A^{-1} b} A^{-1} b d^T A^{-1}$$

$$V_{(n+1)} = V_{(n)} - \frac{1}{1 + \tilde{x}_{n+1}^T V_{(n)} \tilde{x}_{n+1}} V_{(n)} \tilde{x}_{n+1} \tilde{x}_{n+1}^T V_{(n)}$$

$$\hat{\beta}_{(n+1)} = V_{(n+1)} (X_{(n)}^T y + \tilde{x}_{n+1} y_{n+1})$$

$$= \left(V_{(n)} - \frac{1}{1 + \tilde{x}_{n+1}^T V_{(n)} \tilde{x}_{n+1}} V_{(n)} \tilde{x}_{n+1} \tilde{x}_{n+1}^T V_{(n)} \right) (X_{(n)}^T y + \tilde{x}_{n+1} y_{n+1})$$

$$\begin{aligned}
 \hat{\beta}_{(n+1)} &= V_{(n+1)} (X_{(n)}^T y + \tilde{x}_{n+1} y_{n+1}) \\
 &= \left(V_{(n)} - \frac{1}{1 + \tilde{x}_{n+1}^T V_{(n)} \tilde{x}_{n+1}} V_{(n)} \tilde{x}_{n+1} \tilde{x}_{n+1}^T V_{(n)} \right) (X_{(n)}^T y + \tilde{x}_{n+1} y_{n+1}) \\
 &= \underbrace{V_{(n)} X_{(n)}^T y}_{\hat{\beta}_{(n)}} + V_{(n)} \tilde{x}_{n+1} y_{n+1} - h V_{(n)} \tilde{x}_{n+1} \tilde{x}_{n+1}^T V_{(n)} X_{(n)}^T y - h V_{(n)} \tilde{x}_{n+1} \tilde{x}_{n+1}^T V_{(n)} \tilde{x}_{n+1} y_{n+1} \\
 &= \hat{\beta}_{(n)} + h \left(V_{(n)} \tilde{x}_{n+1} y_{n+1} + \tilde{x}_{n+1}^T V_{(n)} \tilde{x}_{n+1} V_{(n)} X_{(n)}^T y - V_{(n)} \tilde{x}_{n+1} \tilde{x}_{n+1}^T \underbrace{V_{(n)} X_{(n)}^T y}_{\hat{\beta}_{(n)}} \right. \\
 &\quad \left. - V_{(n)} \tilde{x}_{n+1} \tilde{x}_{n+1}^T V_{(n)} \tilde{x}_{n+1} y_{n+1} \right) \\
 &= \hat{\beta}_{(n)} + \underbrace{h V_{(n)} \tilde{x}_{n+1}}_{K_n} \underbrace{\left(y_{n+1} - \tilde{x}_{n+1}^T \hat{\beta}_{(n)} \right)}_{e_{n+1}} \\
 &= \hat{\beta}_{(n)} + K_n e_{n+1}
 \end{aligned}$$

\rightarrow prediction error of y_{n+1} based on the estimate obtained on the previous step n ($\hat{\beta}_{(n)}$)

$$\hat{\beta}_{(n+1)} \text{ and } V_{(n+1)} \rightarrow \hat{\beta}_{(n+2)} \text{ and } V_{(n+2)} \rightarrow \dots$$