

# Linear models and misspecification

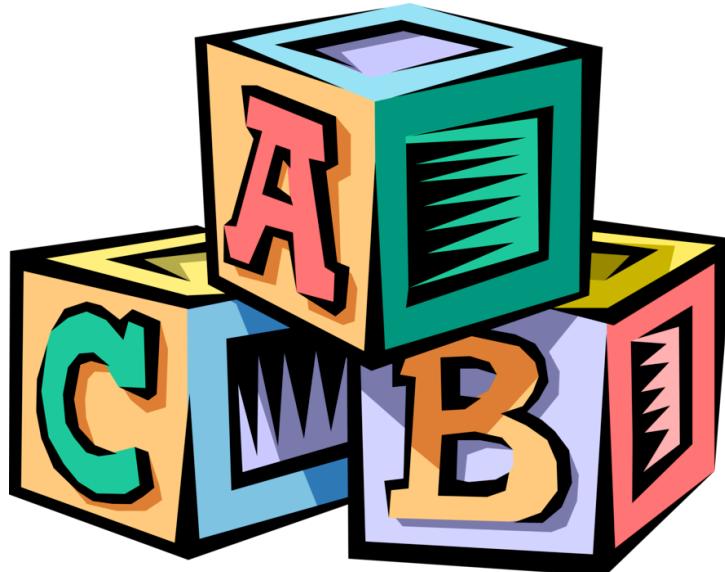
Statistics III - CdL SSE

**Tommaso Rigon**

*Università degli Studi di Milano-Bicocca*

[Home page](#)

# Homepage



*“Everything should be made as simple as possible, but not simpler”*

Attributed to Albert Einstein

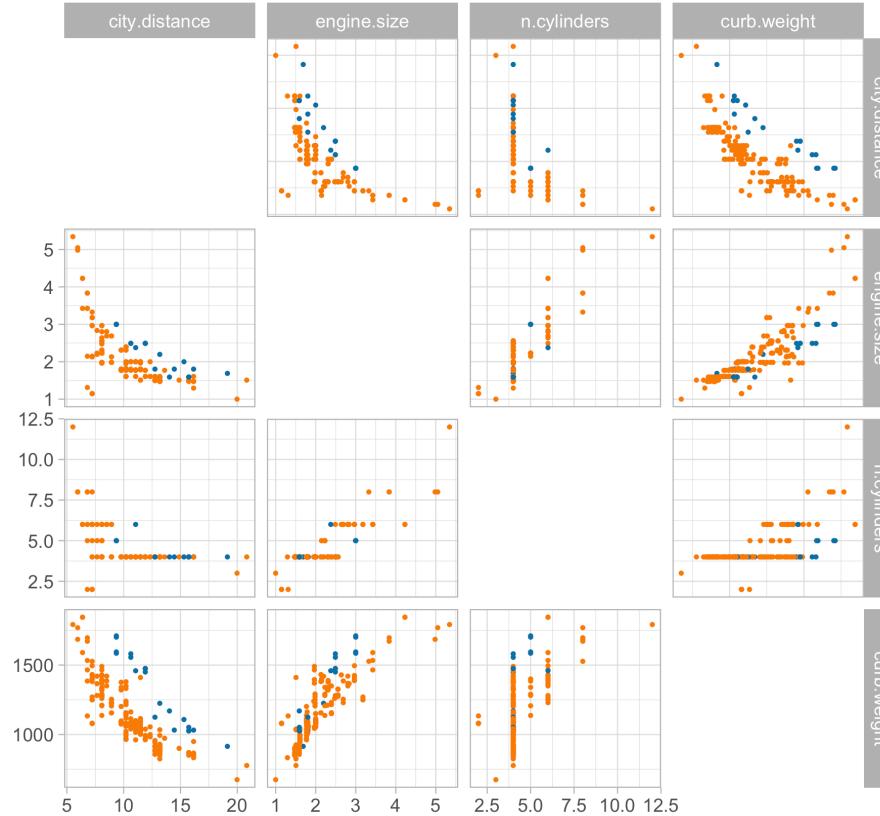
- This unit will cover the following **topics**:
  - **Recap**: linear models and the modeling process
  - Robustness of OLS estimates, sandwich estimators
  - Weighted least squares
  - Box-Cox transform, variance stabilizing transformations
- The main theme is: what should we do when the **assumptions** of linear models are **violated**?
- We will push the linear model to its limit, using it even when it is not supposed to work.
- The symbol means that a few extra steps are discussed in the **handwritten notes**.

The content of this Unit is covered in **Chapter 1** of Salvan et al. (2020). Alternatively, see **Chapter 2** of Agresti (2015) or **Chapter 5** of Azzalini (2008).

# The modeling process

[Home page](#)

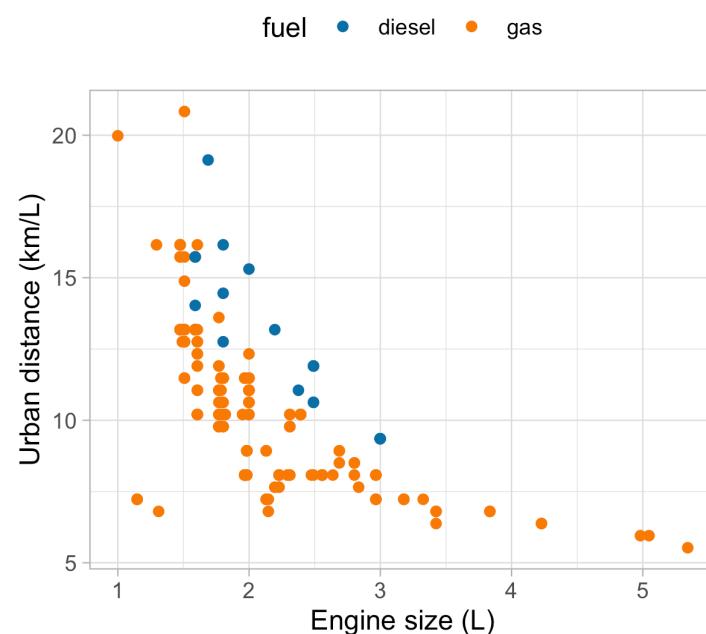
# Car data (diesel or gas)



- We consider data for  $n = 203$  models of cars in circulation in 1985 in the USA.
- We want to **predict** the distance per unit of fuel as a function of the vehicle features.
- We consider the following **variables**:
  - The city distance per unit of fuel (km/L, **city.distance**)
  - The engine size (L, **engine.size**)
  - The number of cylinders (**n.cylinders**)
  - The curb weight (kg, **curb.weight**)
  - The fuel type (gasoline or diesel, **fuel**).

We assume you are already familiar with linear models. The following is a brief recap rather than a full discussion.

# Linear regression



- Let us consider the variables `city.distance` ( $y$ ), `engine.size` ( $x$ ) and `fuel` ( $z$ ).
  - A **simple linear regression**
- $$Y_i = \beta_1 + \beta_2 x_i + \epsilon_i, \quad i = 1, \dots, n,$$
- could be easily fit by least squares...
- ... but the plot suggests that the relationship between `city.distance` and `engine.size` is **not** well approximated by a **linear** function.
  - ... and also that `fuel` has a non-negligible effect on the response.

# Regression models

[Home page](#)

# Linear models

- Let us consider again the variables `city.distance` ( $y$ ), `engine.size` ( $x$ ) and `fuel` ( $z$ ).
- Which function  $f(x, z; \beta)$  should we choose?

# Matrix notation

- The **response random variables** are collected in the random vector  $\mathbf{Y} = (Y_1, \dots, Y_n)^T$ , whose **observed realization** is  $\mathbf{y} = (y_1, \dots, y_n)^T$ .
- The **design matrix** is a  $n \times p$  matrix, comprising the covariate's values, defined by

$$\mathbf{X} = \begin{bmatrix} x_{11} & \cdots & x_{1p} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{np} \end{bmatrix}.$$

# Linear regression: estimation I

- The optimal set of coefficients  $\hat{\beta}$  is the minimizer of the **least squared criterion**

$$D(\beta) = (\mathbf{y} - \mathbf{X}\beta)^T(\mathbf{y} - \mathbf{X}\beta) = \|\mathbf{y} - \mathbf{X}\beta\|^2,$$

also known as **residual sum of squares (RSS)**, where

$$\|\mathbf{y}\| = \sqrt{y_1^2 + \cdots + y_n^2},$$

denotes the **Euclidean norm**.

## Linear regression: estimation II

- In matrix notation, the predicted values can be obtained as

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\beta} = \mathbf{H}\mathbf{y}, \quad \mathbf{H} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T,$$

where  $\mathbf{H}$  is a  $n \times n$  **projection matrix** matrix sometimes called **hat matrix**. The matrix is idempotent, meaning that  $\mathbf{H} = \mathbf{H}^T$  and  $\mathbf{H}^2 = \mathbf{H}$ .

## Linear regression: inference

- Recall that the errors  $\epsilon$  have zero mean  $\mathbb{E}(\epsilon) = 0$  and are **uncorrelated**  $\text{var}(\epsilon) = \sigma^2 I_n$ .
- Then, the estimator  $\hat{\beta}$  is **unbiased**  $\mathbb{E}(\hat{\beta}) = \beta$  and its **variance** is  $\text{var}(\hat{\beta}) = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$ . Since  $\sigma^2$  is also unknown, we can estimate the variances of  $\hat{\beta}$  as follows:

$$\widehat{\text{var}}(\hat{\beta}) = s^2 (\mathbf{X}^T \mathbf{X})^{-1}.$$

- The **standard errors** of the components of  $\hat{\beta}$  correspond to the square root of the diagonal of the above covariance matrix.

# Linear regression: diagnostic

- The diagonal elements  $h_i \in [0, 1]$  of the matrix  $\mathbf{H}$  are called **leverages** and it holds

$$\text{var}(\hat{Y}_i) = \sigma^2 h_i, \quad \text{var}(Y_i - \hat{Y}_i) = \sigma^2(1 - h_i), \quad \text{cor}(Y_i, \hat{Y}_i) = \sqrt{h_i}.$$

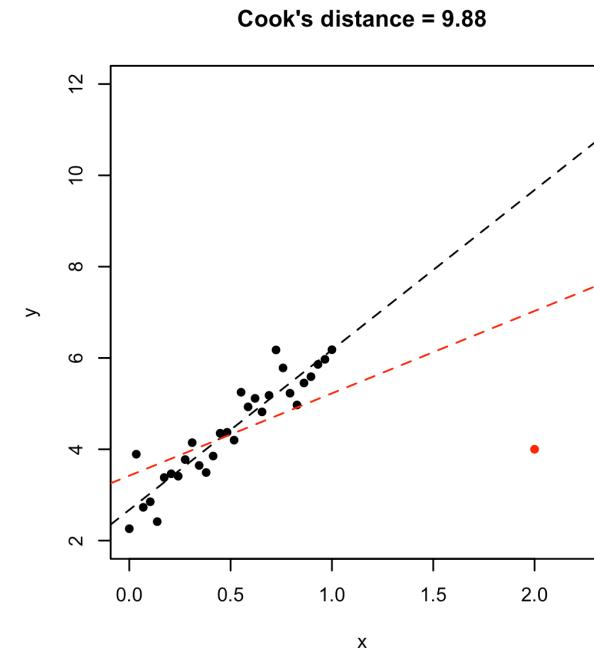
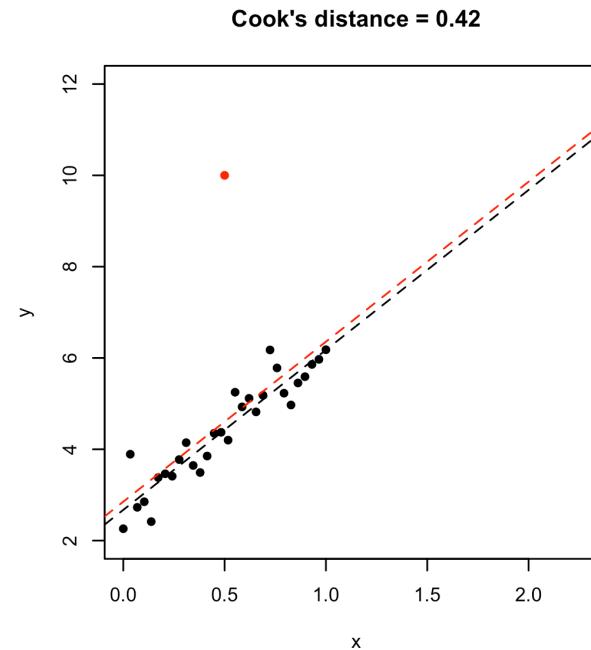
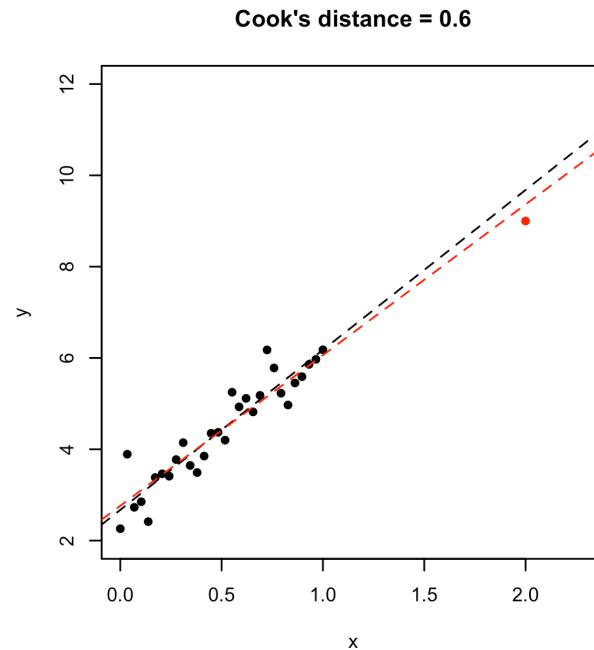
The leverage  $h_i$  determines the **precision** with which  $\hat{Y}_i$  predicts  $Y_i$ . For large  $h_i$  close to 1,  $\text{cor}(Y_i, \hat{Y}_i) \approx 1$ , therefore changes of a single point  $Y_i$  leads to significant changes in  $\hat{Y}_i$ .

- Leverages also appear in the definition of **standardized residuals**:

$$\tilde{r}_i = \frac{r_i}{\sqrt{s^2(1 - h_i)}} = \frac{y_i - \mathbf{x}_i^T \hat{\beta}}{\sqrt{s^2(1 - h_i)}},$$

where  $r_i = y_i - \mathbf{x}_i^T \hat{\beta}$  are the (raw) **residuals**.

## Leverages, outliers and influence points



- **Left plot:** leverage, not outlier. **Central plot:** outlier, not leverage. **Right plot:** influence point = leverage + outlier.

## A first model: estimated coefficients

- Our first attempt for predicting `city.distance` ( $y$ ) via `engine.size` ( $x$ ) and `fuel` ( $z$ ) is:

$$f(x, z; \beta) = \beta_1 + \beta_2 x + \beta_3 x^2 + \beta_4 x^3 + \beta_5 I(z = \text{gas}).$$

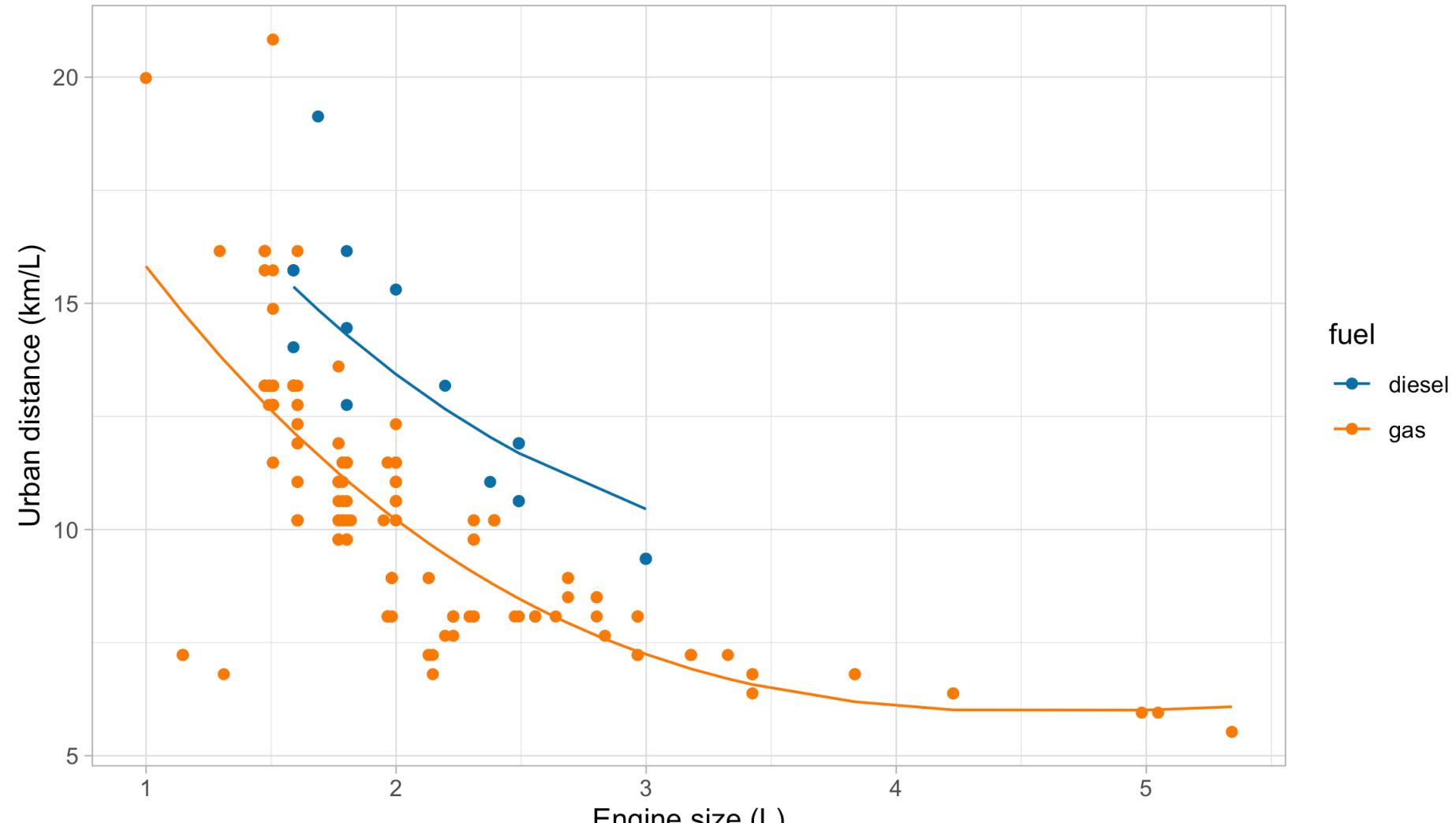
- We obtain the following **summary** for the regression coefficients  $\hat{\beta}$ .

term	estimate	std.error	statistic	p.value
(Intercept)	28.045	3.076	9.119	0.000
engine.size	-10.980	3.531	-3.109	0.002
engine.size^2	2.098	1.271	1.651	0.100
engine.size^3	-0.131	0.139	-0.939	0.349
fuel_gas	-3.214	0.427	-7.523	0.000

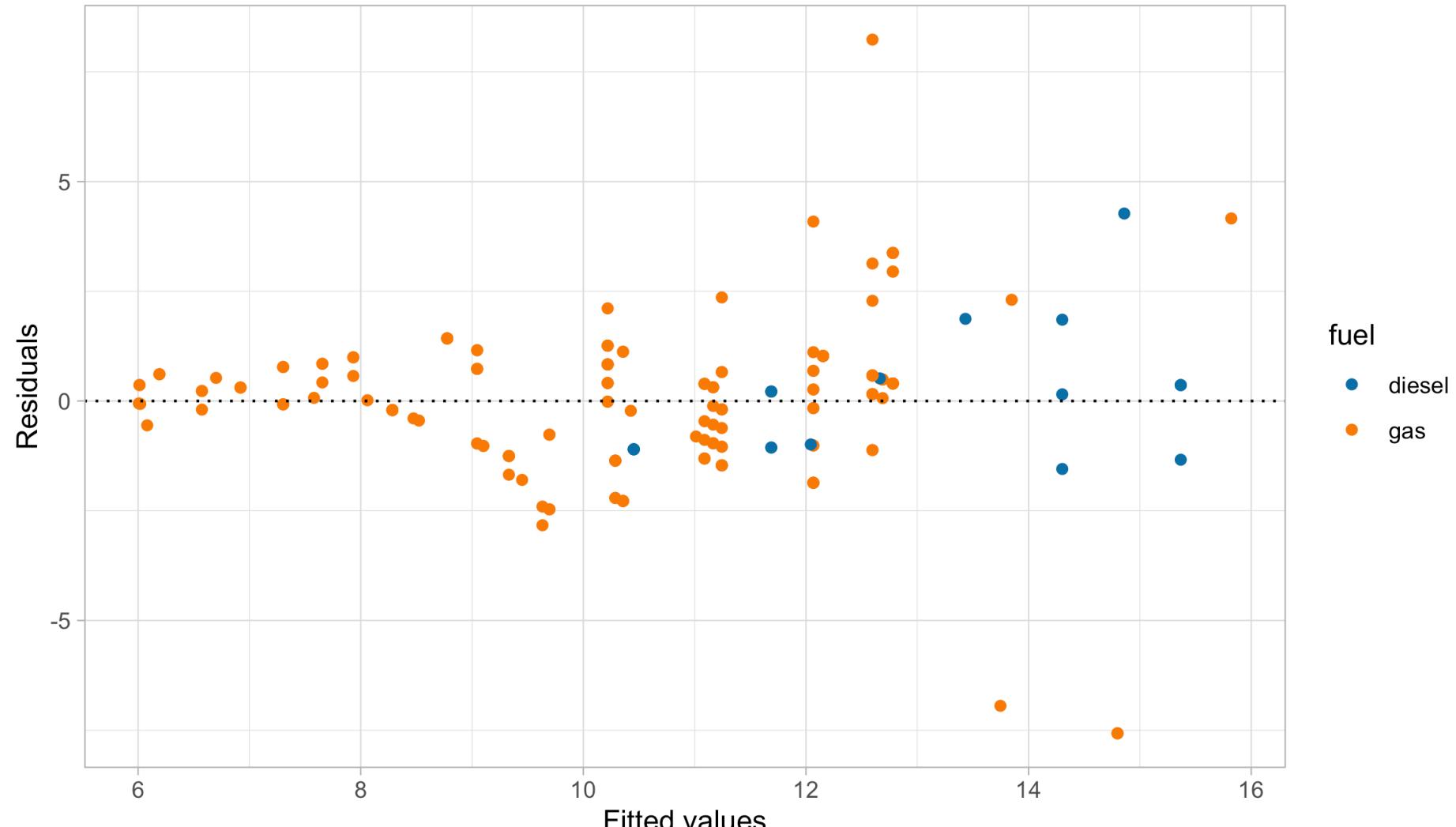
- Moreover, the coefficient  $R^2$  and the residual standard deviation  $s$  are:

r.squared	sigma	deviance
0.5973454	1.790362	634.6687

## A first model: fitted values



# A first model: graphical diagnostics



Home page

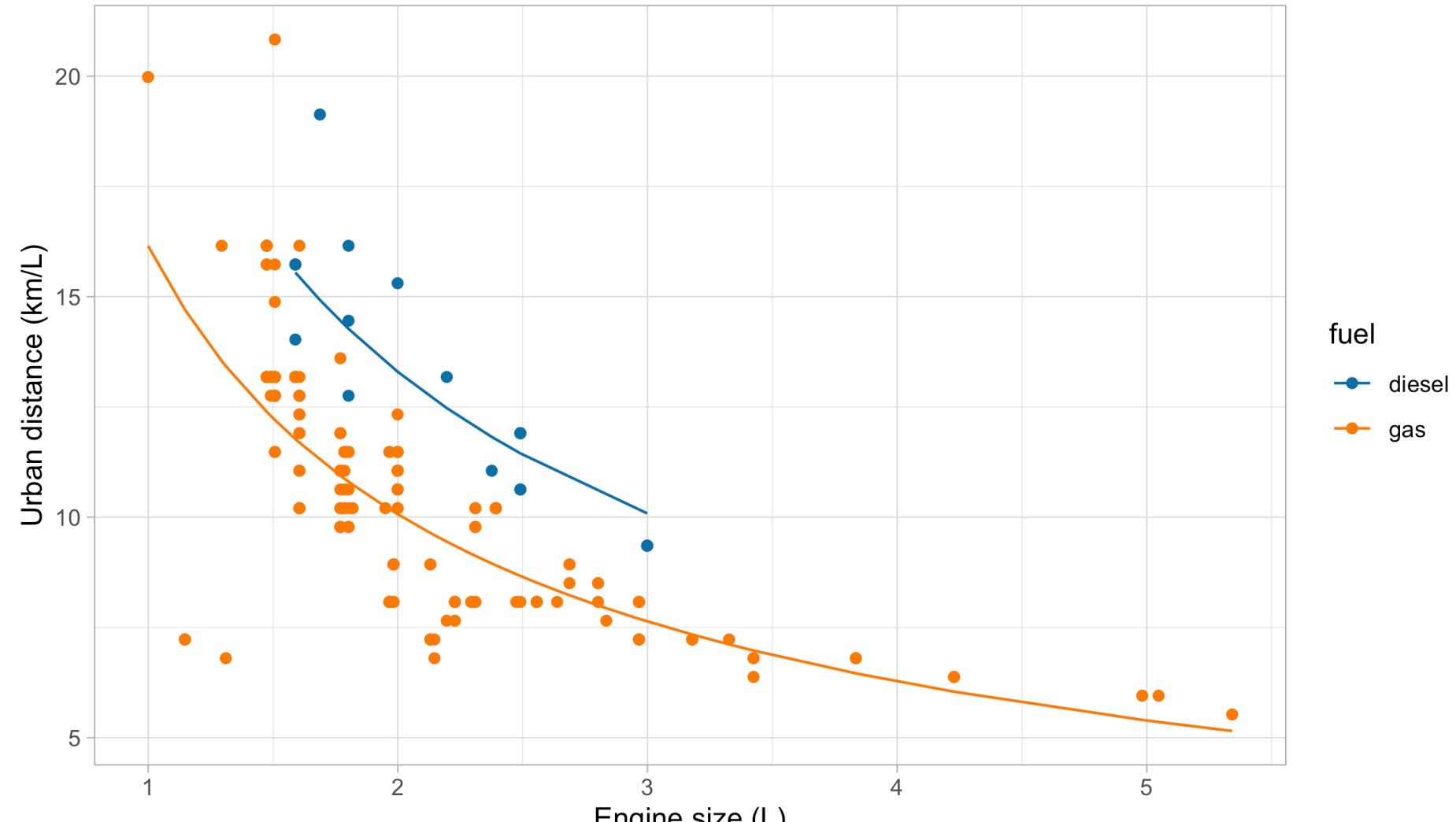
## Comments and criticisms

- Is this a good model?
- The overall fit **seems satisfactory** at first glance, especially if we aim at predicting the urban distance of cars when average engine size (i.e., between  $1.5L$  and  $3L$ ).

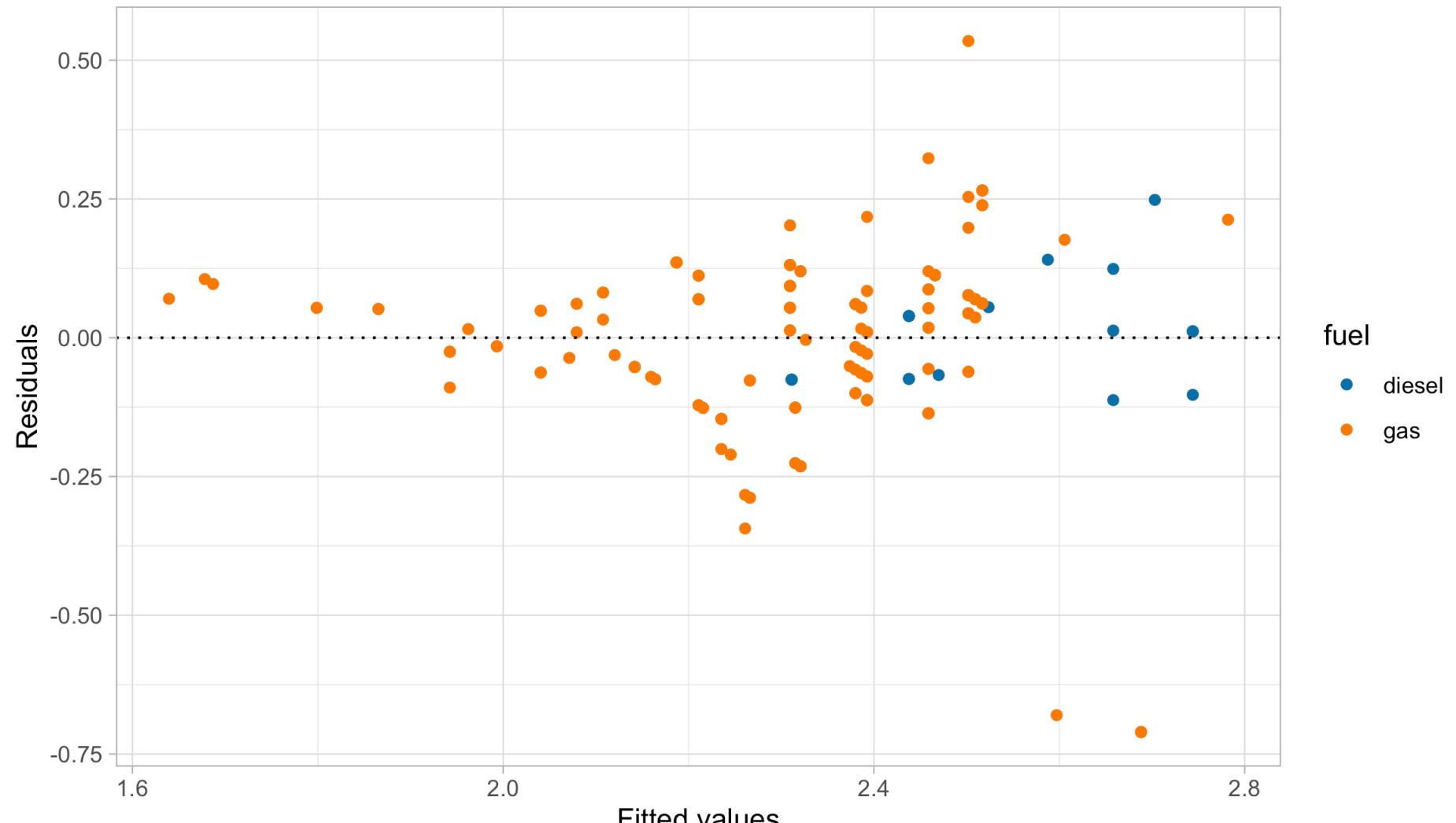
## Linear models and non-linear patterns

- A significant advantage of linear models is that they can describe non-linear relationships via **variable transformations** such as polynomials, logarithms, etc.

## Second model: fitted values



## Second model: graphical diagnostics



## Comments and criticisms

- The **goodness of fit** indices are the following:

r.squared.original	r.squared	sigma	deviance
0.5847555	0.6196093	0.1600278	5.121777

- Do not mix **apple** and **oranges**! Compare  $R^2$ s only if they refer to the same scale!

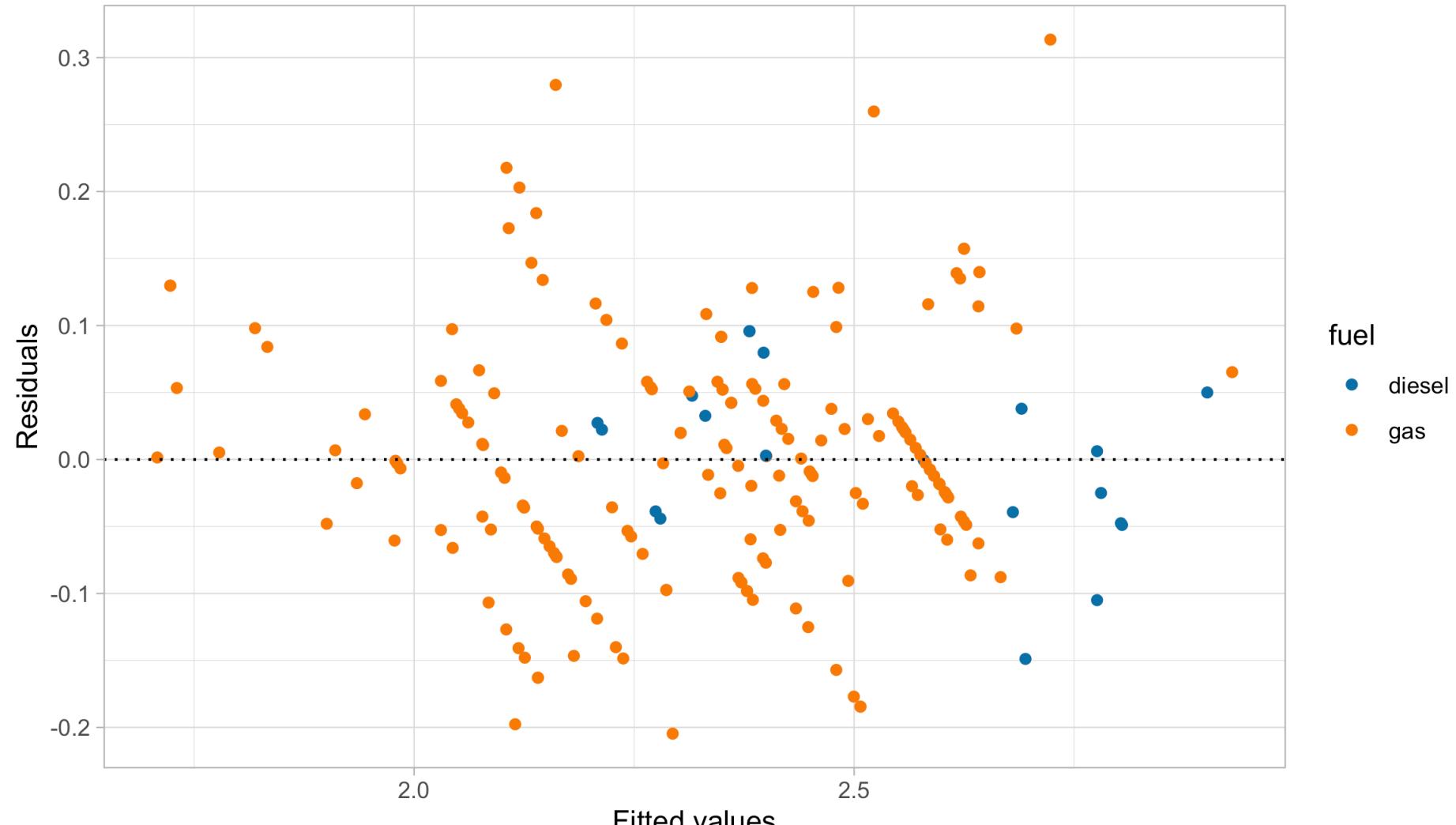
## A third model: additional variables

- Let us consider **two additional variables**: `curb.weight` ( $w$ ) and `n.cylinders` ( $v$ ).
- A richer model, therefore, could be:

$$\log Y_i = \beta_1 + \beta_2 \log x_i + \beta_3 \log w_i + \beta_4 I(z_i = \text{gas}) + \beta_5 I(v_i = 2) + \epsilon_i,$$

for  $i = 1, \dots, n$ . The estimates are:

## A third model: graphical diagnostics



## Comments and criticisms

- The goodness of fit greatly **improved**:

r.squared.original	r.squared	sigma	deviance
0.869048	0.8819199	0.0896089	1.589891

- In this third model, we handled the **outliers** appearing in the residual plots, which it turns out are identified by the group of cars having 2 cylinders.
- The diagnostic plots are also very much improved, although still not perfect.
- The estimates are coherent with our expectations, based on common knowledge. Have a look at the book (Azzalini and Scarpa (2012)) for a detailed explanation of  $\beta_4$ !
- The car dataset is available from the textbook (A&S) website:
  - Dataset <http://azzalini.stat.unipd.it/Book-DM/auto.dat>
  - Variable description <http://azzalini.stat.unipd.it/Book-DM/auto.names>

# Misspecification and remedies

[Home page](#)

# Assumptions and misspecification

## Classical assumptions of linear models

- (A.1) **Linear structure**, namely  $\mathbf{Y} = \mathbf{X}\beta + \boldsymbol{\epsilon}$  with  $\mathbb{E}(\boldsymbol{\epsilon}) = 0$ , implying  $\mathbb{E}(\mathbf{Y}) = \mathbf{X}\beta$ .<sup>1</sup>
- (A.2) **Homoschedasticity** and **uncorrelation** of the errors, namely  $\text{var}(\boldsymbol{\epsilon}) = \sigma^2 I_n$ .
- (A.3) **Gaussianity**, namely  $\boldsymbol{\epsilon} \sim N_n(0, \sigma^2 I_n)$ . In other words, the errors  $\epsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$  are iid Gaussian random variables with zero mean and variance  $\sigma^2$ .

It is also commonly asked that  $\text{rk}(\mathbf{X}) = p$ , otherwise the model is not identifiable.

- If one of the above assumptions is violated, it is not necessarily a huge problem, because
    - the OLS estimator  $\hat{\beta}$  is fairly **robust** to misspecification;
    - simple **fixes** (variable transformations, standard error corrections) are available.
1. If the intercept is included in  $\mathbf{X}$ , the errors automatically satisfy the property  $\mathbb{E}(\boldsymbol{\epsilon}) = 0$ .

# Robust estimation and assumptions



- A plane can still fly with one of its **engines on fire**, but this is hardly an appealing situation.
- Similarly, robust estimators may work under **model misspecification**, but this does not mean we should neglect **checking** whether the original **assumptions** hold.

# Non-normality of the errors I



- Let us consider the case in which assumptions **(A.1)-(A.2)** are **valid** but **(A.3) is not**, that is  $\mathbb{E}(\epsilon) = 0$  and  $\text{var}(\epsilon) = \sigma^2 I_n$ , but  $\epsilon$  does **not** follow **a Gaussian** distribution.
- For example,  $\epsilon_i$  may follow a Laplace distribution, a skew-Normal, a logistic distribution, a Student's t distribution, etc.

## Non-normality of the errors II

- When the errors are non Gaussian the **exact inferential results** are not valid. In particular  $\hat{\beta}$  does not follow anymore a Gaussian distribution.
- However, a **central limit theorem** can be invoked under very mild conditions on the design matrix  $\mathbf{X}$ .
- Thus, when the sample size  $n$  is large enough, then the following **approximation** holds

$$\hat{\beta} \stackrel{\text{d}}{\sim} N_p(\beta, \sigma^2(\mathbf{X}^T \mathbf{X})^{-1}),$$

from which **confidence intervals** and **test statistics** can be obtained as usual. The approximation is **excellent** if the errors are **symmetric** around 0.

# Heteroschedasticity of the errors I



- Suppose now that the linearity assumption **(A.1)** is valid but **homoschedasticity** of the errors **(A.2)** is **not**. Instead, we consider **heteroschedastic errors**:

$$\text{var}(\epsilon) = \Sigma, \quad \text{or equivalently} \quad \text{var}(Y_i) = \sigma_i^2, \quad i = 1, \dots, n$$

where  $\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_n^2)$  is a diagonal matrix with positive entries.

- The OLS estimator is still **unbiased**, with a **modified covariance** structure<sup>1</sup>

$$\mathbb{E}(\hat{\beta}) = \beta, \quad \text{var}(\hat{\beta}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \Sigma \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}.$$

If in addition we assume Gaussianity of the errors, that is  $\epsilon \sim N_n(0, \Sigma)$ , then

$$\hat{\beta} \sim N_p(\beta, (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \Sigma \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}).$$

Under suitable but mild conditions on  $\mathbf{X}$  and  $\Sigma$ , the estimator is also **consistent**.

- These results are valid even when the matrix  $\Sigma$  is non-diagonal. This is useful to model correlated responses.

## Heteroschedasticity of the errors II

The OLS estimator in presence of heteroschedasticity still gives a **good point estimate**. However, the OLS estimator is **not efficient** and the classical **standard errors** are **wrong**.

- A potential approach is to **accept the inefficiency** of the OLS estimator in this scenario and **correct** the standard errors.
- The elements of  $\Sigma$  are **unknown**, but we can estimate them from the data. Note that

$$\text{var}(r_i) = \text{var}(y_i - \mathbf{x}_i^T \hat{\beta}) = \sigma_i^2(1 - h_i),$$

suggesting the **estimate**  $\hat{\sigma}_i^2 = r_i^2/(1 - h_i)$ .

- This leads to the so-called **sandwich estimator** of the covariance matrix:

$$\widehat{\text{var}}(\hat{\beta}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \hat{\Sigma} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1},$$

where  $\hat{\Sigma} = \text{diag}(\hat{w}_1, \dots, \hat{w}_n)$  and  $\hat{w}_i = r_i^2/(1 - h_i)$ .

- These are known as **White's** heteroscedasticity-consistent **standard errors**.<sup>1</sup>

1. White originally proposed the simpler version  $\hat{\sigma}_i^2 = r_i^2$ . Another variant is  $\hat{\sigma}_i^2 = r_i^2/(1 - h_i)^2$ .

# Weighted least squares I



- Let us consider again the case of **heteroschedastic errors**:

$$\text{var}(\epsilon) = \sigma^2 \Omega^{-1}, \quad \text{or equivalently} \quad \text{var}(Y_i) = \sigma_i^2 = \frac{\sigma^2}{\omega_i}, \quad i = 1, \dots, n$$

where  $\Omega = \text{diag}(\omega_1, \dots, \omega_n)$  are positive **weights**. However, here we assume that the weights  $\omega_1, \dots, \omega_n$  are **known**, a common situation in survey design.

- Let us define the **standardized** quantities:

$$\mathbf{Y}^* = \Omega^{1/2} \mathbf{Y}, \quad \mathbf{X}^* = \Omega^{1/2} \mathbf{X}.$$

This is equivalent to say that  $Y_i^* = \sqrt{\omega_i} Y_i$  and  $x_{ij}^* = \sqrt{\omega_i} x_{ij}$ . Then, it is easy to show that

$$\mathbb{E}(\mathbf{Y}^*) = \mathbf{X}^* \beta, \quad \text{var}(\mathbf{Y}^*) = \sigma^2 \Omega^{1/2} \Omega^{-1} \Omega^{1/2} = \sigma^2 I_n,$$

namely the **assumptions (A.1)** and **(A.2)** are valid in the **transformed scale**.

- In other words, **after** a suitable **transformation**, we reconducted the problem to a **standard linear model**.

## Weighted least squares II

- Thus an estimator for  $\beta$ , based on the transformed data, is obtained minimizing the deviance

$$\begin{aligned} D_{\text{wls}}(\beta) &= (\mathbf{y}^* - \mathbf{X}^* \beta)^T (\mathbf{y}^* - \mathbf{X}^* \beta) = (\mathbf{y} - \mathbf{X} \beta)^T \boldsymbol{\Omega} (\mathbf{y} - \mathbf{X} \beta) \\ &= \sum_{i=1}^n \omega_i (y_i - \mathbf{x}_i^T \beta)^2. \end{aligned}$$

which is a **weighted** version of the original **quadratic loss**, with **high weight = low variance**.

- The resulting OLS estimate minimizing  $D_{\text{wls}}(\beta)$  in the transformed and original scales is

$$\hat{\beta}_{\text{wls}} = [(\mathbf{X}^*)^T \mathbf{X}^*]^{-1} (\mathbf{X}^*)^T \mathbf{y}^* = (\mathbf{X}^T \boldsymbol{\Omega} \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\Omega} \mathbf{y}$$

and it is referred to as **weighted least squares** estimator of  $\beta$ .

- Such an estimator is **unbiased** and **efficient (BLUE)**, with

$$\mathbb{E}(\hat{\beta}_{\text{wls}}) = \beta, \quad \text{var}(\hat{\beta}_{\text{wls}}) = \sigma^2 (\mathbf{X}^T \boldsymbol{\Omega} \mathbf{X})^{-1}.$$

Moreover, if  $\epsilon \sim N_n(0, \sigma^2 \boldsymbol{\Omega}^{-1})$  it also coincides with the **maximum likelihood** estimator.

# Variable transformations

- Another remedy for **misspecification** was already applied in the analysis of the car dataset, namely through **variable transformation**.

While the model may have been incorrectly specified for the original data, it could become **appropriate** once the **transformations** are considered, namely

$$g(Y_i) = h_1(\mathbf{x}_i)\beta_1 + \cdots + h_p(\mathbf{x}_i)\beta_p + \epsilon_i, \quad i = 1, \dots, n,$$

where  $g(\cdot)$  and  $h_j(\cdot)$  for  $j = 1, \dots, p$  are **non-linear** and **known** functions.

# Box-Cox transform

## Box-Cox transform

If the data are  $y_i$  are **positive**, we may consider a **parametric class** of transformations:

$$g_\lambda(y) = \frac{y^\lambda - 1}{\lambda}, \quad \lambda \neq 0.$$

and  $g_\lambda(y) = \log y$  when  $\lambda = 0$ . This is the celebrated **Box-Cox transform**.

The case  $\lambda = 1$  corresponds to no transformation,  $\lambda = 1/2$  to the square root,  $\lambda = 0$  to the logarithm, and  $\lambda = -1$  to the reciprocal.

## Box-Cox transform: derivation I

- By assumption, the distribution of the **transformed data**  $\mathbf{Z}_\lambda = (g_\lambda(Y_1), \dots, g_\lambda(Y_n))^T$  is Gaussian, therefore their joint density is

$$f_Z(\mathbf{z}_\lambda) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp \left\{ -\frac{1}{2\sigma^2} (\mathbf{z}_\lambda - \mathbf{X}\beta)^T (\mathbf{z}_\lambda - \mathbf{X}\beta) \right\}.$$

- Using standard tools of probability theory, we can obtain the density of the **original data**:

$$f_Y(\mathbf{y}) = f_Z(g_\lambda(y_1), \dots, g_\lambda(y_n)) \prod_{i=1}^n \left| \frac{\partial g_\lambda(y_i)}{\partial y_i} \right|, \quad \text{where} \quad \left| \frac{\partial g_\lambda(y_i)}{\partial y_i} \right| = y_i^{\lambda-1}.$$

The additional term is the determinant of the **Jacobian** of the transformation.

- The **log-likelihood** therefore is

$$\ell(\beta, \sigma^2, \lambda) = -\frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (\mathbf{z}_\lambda - \mathbf{X}\beta)^T (\mathbf{z}_\lambda - \mathbf{X}\beta) + (\lambda - 1) \sum_{i=1}^n \log y_i.$$

## Box-Cox transform: derivation II

- Note that, for any given value of  $\lambda$ , the maximum likelihood estimates are

$$\hat{\beta}_\lambda = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{z}_\lambda, \quad \hat{\sigma}_\lambda^2 = \frac{1}{n} (\mathbf{z}_\lambda - \mathbf{X} \hat{\beta}_\lambda)^T (\mathbf{z}_\lambda - \mathbf{X} \hat{\beta}_\lambda),$$

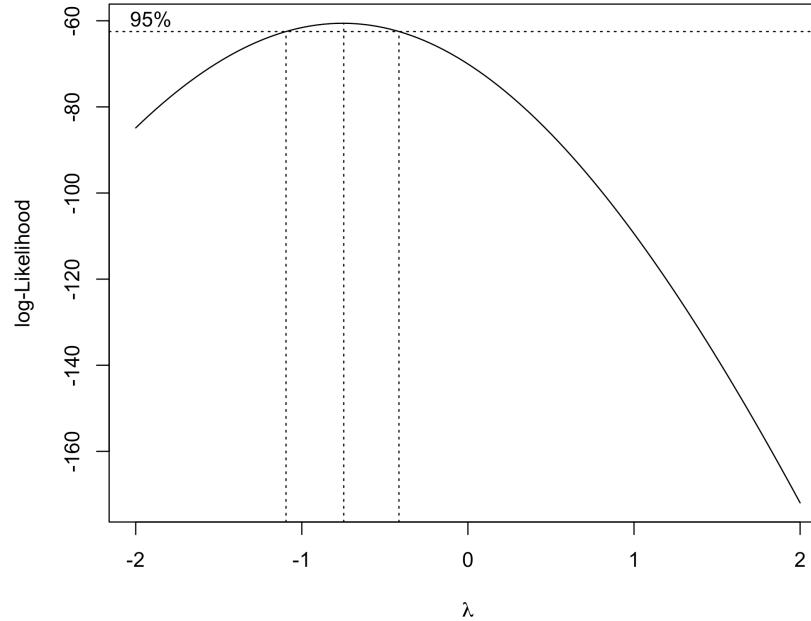
- We can **plug-in** the above estimates into the log-likelihood. This gives the **profile log-likelihood** for  $\lambda$ , which admits a very simple expression:

$$\ell_P(\lambda) = \ell(\hat{\beta}_\lambda, \hat{\sigma}_\lambda^2, \lambda) = -\frac{n}{2} \log \hat{\sigma}_\lambda^2 + (\lambda - 1) \sum_{i=1}^n \log y_i,$$

which must be **numerically maximized** over  $\lambda$ , e.g. using `optim`.

- The optimal value  $\hat{\lambda} = \arg \max \ell_P(\lambda)$ , as well as a confidence interval for it, may offer guidance in choosing the right transformation.

# Box-Cox transform for the auto dataset

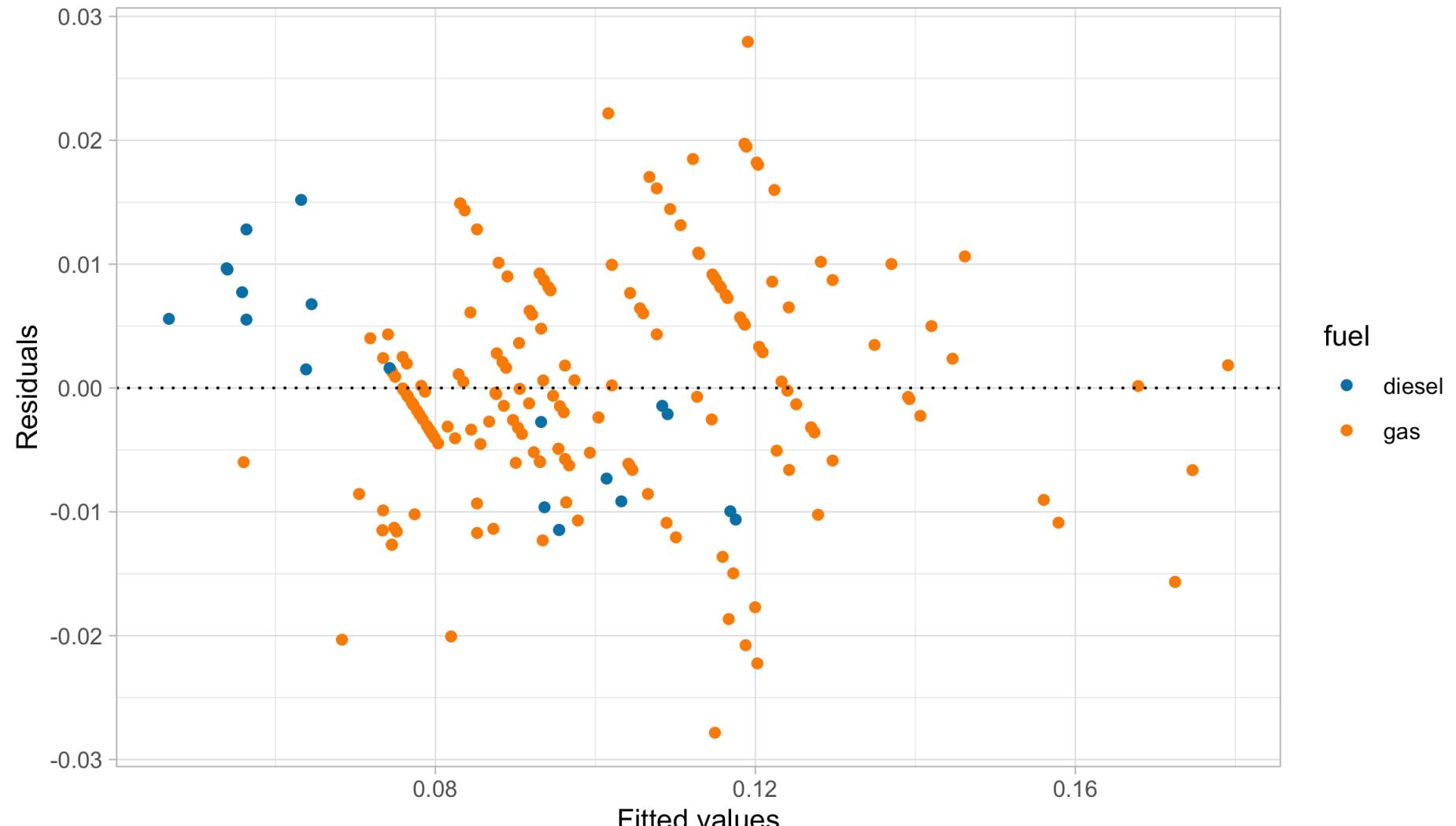


- The **Box-Cox transform** in the auto dataset suggests a **reciprocal** transformation:

$$\frac{1}{Y_i} = \beta_1 + \beta_2 x_i + \beta_3 w_i + \beta_4 I(z_i = \text{gas}) + \beta_5 I(v_i = 2) + \epsilon_i,$$

which is a good alternative to our model based on logarithms of  $y_i$ ,  $x_i$ , and  $w_i$  (**but...**).

## A fourth model: graphical diagnostics



# Variance stabilizing transformations I



- Let  $Y_i \sim \text{Poisson}(\mu_i)$  with mean  $\mathbb{E}(Y_i) = \mu_i = f(\mathbf{x}_i; \beta) = \text{var}(Y_i)$ . Note that

$$Y_i \stackrel{\text{d}}{\sim} \text{N}(\mu_i, \mu_i),$$

is **asymptotically Gaussian** for large values of  $\mu_i$ . However, data are **heteroschedastic**.

- In modeling count data, we could transform the counts so that, at least **approximately**, the **variance** of  $g(Y_i)$  is **constant** and ordinary least squares methods can be used.

# Variance stabilizing transformations II



- Let  $Y_i \sim \text{Binomial}(\pi_i, m_i)$ , with **success probability**  $\pi_i = f(x_i; \beta)$  and **trials**  $m_i$ . For large values of  $m_i$ , the **Gaussian approximation** holds

$$Y_i \stackrel{\text{d}}{\sim} N(m_i\pi_i, m_i\pi_i(1 - \pi_i)).$$

However, the data are **heteroschedastic**, because  $\text{var}(Y_i) = m_i\pi_i(1 - \pi_i)$ .

- Thus, a **variance stabilizing** transformation in this case is

$$g_{m_i}(y) = \sqrt{m_i} \arcsin \left( \frac{2y}{m_i} - 1 \right),$$

because in fact we have that

$$\text{var}(g_{m_i}(Y_i)) \approx \left( \frac{\sqrt{m_i}}{\sqrt{1 - (2\pi_i - 1)^2}} \frac{2}{m_i} \right)^2 m_i\pi_i(1 - \pi_i) = 1.$$

# Limitations of variable transformations I

- Variable transformations are appealing for their simplicity and have a long history in statistics. However, they also have some **drawbacks**.

## Limitations of variable transformations II

Suppose  $Y_i \sim \text{Binomial}(\pi, m_i)$ . The variance stabilizing transformation is not fully satisfactory:

- It **complicates** the **interpretation**, because it models  $\mathbb{E}\{g(Y_i)\}$  instead of  $\mathbb{E}(Y_i)$ ;
- It is an **asymptotic approximation**, and is only valid for  $m_i \rightarrow \infty$ .
- The transform depends on  $m_i$ , therefore we cannot make predictions for a generic covariate value  $x_i$  without knowing the associated  $m_i$ .

Besides, this transform is clearly not applicable when  $m_i = 1$  and  $Y_i \in \{0, 1\}$ , a very common problem called **binary regression**.

- If we know that  $Y_i$  follows, say, a Bernoulli or a Gamma distribution, then we should use the **appropriate likelihood** rather than a **Gaussian approximation**.
- **Generalized Linear Models** provide a **much more elegant solution** to the above problem.

## References

- Agresti, A. (2015), *Foundations of Linear and Generalized Linear Models*, Wiley.
- Azzalini, A. (2008), *Inferenza statistica*, Springer Verlag.
- Azzalini, A., and Scarpa, B. (2012), *Data analysis and data mining: An introduction*, Oxford University Press.
- Salvan, A., Sartori, N., and Pace, L. (2020), *Modelli lineari generalizzati*, Springer.