

MODEL CHECKING / DIAGNOSTICS

The inference we obtained was performed under the assumption that the hypotheses were met. However, we have to make sure it is the case.

After we fit the model, we need to evaluate the validity of the model.

We should assess whether the model satisfies the underlying assumptions:

1. normality $Y_i \sim N(\mu_i, \sigma^2) \quad i=1, \dots, n$
2. linearity $\mu_i = \beta_0 + \beta_1 x_i$
3. homoscedasticity $\text{var}(Y_i) = \sigma^2$ for all $i=1, \dots, n$
4. independence $\text{cov}(Y_i, Y_k) = 0$ for $i \neq k$

or, equivalently, $Y_i = \mu_i + \epsilon_i \quad i=1, \dots, n$; and we formulate the hypotheses on the errors ϵ_i ($\epsilon_i \stackrel{i.i.d.}{\sim} N(0, \sigma^2) \quad i=1, \dots, n$).

Other possible issues to evaluate are:

- is the functional form adequate? The model may be missing needed covariates, or nonlinear transformations of the variables
- is there any outlier? Unusual observations may have too much influence on the model fit.

(we will focus more on these issues with the exercises)

Classical tools to perform the model's diagnostics:

- visual inspection (plots)
- tests

ANALYSIS OF RESIDUALS

We make assumptions on the model's error terms ϵ_i , which are not observable.

However, after we estimate the model, we can compute the RESIDUALS, which are the "analogous" sample quantity (NOT an estimate!).

The assumptions on ϵ_i have implications on the properties of e_i :

\Rightarrow if the properties of ϵ_i do not hold for the estimated model, we conclude that the hypotheses on ϵ_i were not satisfied.

The residuals are $e_i = y_i - \hat{y}_i \quad i=1, \dots, n$.

We have already showed some properties of e_i :

- a) zero mean: $\frac{1}{n} \sum_{i=1}^n e_i = 0$ (see Lecture 3)
- b) orthogonality w.r.t x : $\sum_{i=1}^n x_i e_i = 0$
indeed, $\sum_{i=1}^n x_i e_i = \sum_{i=1}^n x_i (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) \leftarrow$ 2nd likelihood equation
- c) orthogonality w.r.t \hat{y} : $\sum_{i=1}^n e_i \hat{y}_i = 0$
indeed, $\sum_{i=1}^n e_i \hat{y}_i = \sum_{i=1}^n e_i (\hat{\beta}_0 + \hat{\beta}_1 x_i) = \hat{\beta}_0 \sum_{i=1}^n e_i + \hat{\beta}_1 \sum_{i=1}^n e_i x_i = 0$ (a) (b)
- d) $\text{cov}(x, e) = 0$
indeed, $\text{cov}(x, e) = 0 \Leftrightarrow \text{cov}(x, \epsilon) = 0$
 $\text{cov}(x, \epsilon) = \sum_{i=1}^n (x_i - \bar{x})(\epsilon_i - \bar{\epsilon}) = \sum_{i=1}^n x_i \epsilon_i - \bar{x} \sum_{i=1}^n \epsilon_i = 0$ (b) (a)

Before observing the data, we have the random variables $\epsilon_i = Y_i - \hat{Y}_i \quad i=1, \dots, n$.

DISTRIBUTION of ϵ_i

i. they have normal distribution

$$\epsilon_i = Y_i - \hat{Y}_i = Y_i - \hat{\beta}_0(Y) - \hat{\beta}_1(Y) x_i = Y_i - \sum_{k=1}^n \gamma_k Y_k - x_i \sum_{k=1}^n \omega_k Y_k = \sum_{k=1}^n c_k Y_k$$

for some constants c_k .

Hence ϵ_i is a linear combination of normal r.v.'s $\Rightarrow \epsilon_i \sim N(\cdot, \cdot)$ normal

ii. $E[\epsilon_i] = 0$

$$E[\epsilon_i] = E[Y_i - \hat{Y}_i] = E[Y_i] - E[\hat{Y}_i] = \beta_0 + \beta_1 x_i - E[\hat{\beta}_0 + \hat{\beta}_1 x_i] = \beta_0 + \beta_1 x_i - \beta_0 - \beta_1 x_i = 0$$

iii. $\text{var}(\epsilon_i) = \sigma^2(1 - h_i)$

$$\text{with } h_i = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum_{k=1}^n (x_k - \bar{x})^2} \quad h_i \text{ is called "LEVERAGES"}$$

\Rightarrow NOT homoscedastic! (they depend on the index i)

Moreover, they are NOT independent

• Distribution of the residuals: $\epsilon_i \sim N(0, \sigma^2(1 - h_i)) \quad i=1, \dots, n$

ALTERNATIVE DEFINITIONS:

- Standardized residuals $\tilde{\epsilon}_i = \frac{\epsilon_i}{\sqrt{1 - h_i}}$ with $E[\tilde{\epsilon}_i] = 0, \text{var}(\tilde{\epsilon}_i) = \sigma^2$
 $\tilde{\epsilon}_i \sim N(0, \sigma^2)$
homoscedastic, but σ^2 is unknown
- Studentized residuals $R_i = \frac{\epsilon_i}{\hat{\sigma} \sqrt{1 - h_i}}$ with $E[R_i] = 0, \text{var}(R_i) = 1$
we don't have a nice exact distribution, but approximately $R_i \sim N(0, 1)$

How can we use these quantities to check the model's assumptions?

(1) ORTHOGONALITY w.r.t x and HOMOSEDASTICITY

RESIDUALS vs COVARIATE \rightarrow SCATTERPLOT of x_i vs $\tilde{\epsilon}_i$ (standardized residuals)

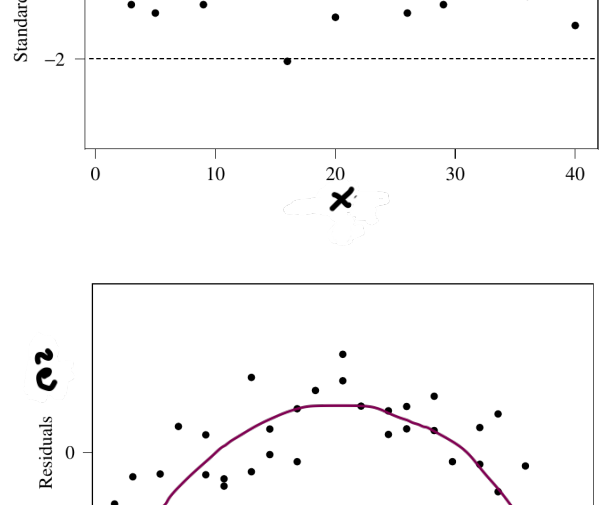
we can use it to check

- the covariate and residuals have correlation = 0
- the standardized residuals are homoscedastic

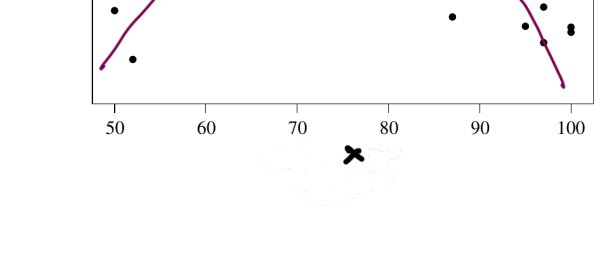
note: we use standardized or studentized residuals to have constant variance

if the assumption is satisfied, the plot should show a random pattern

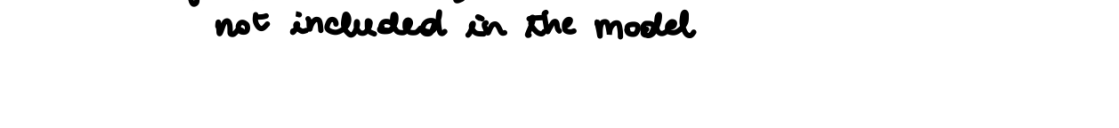
(no systematic behaviors) and homogeneous variability



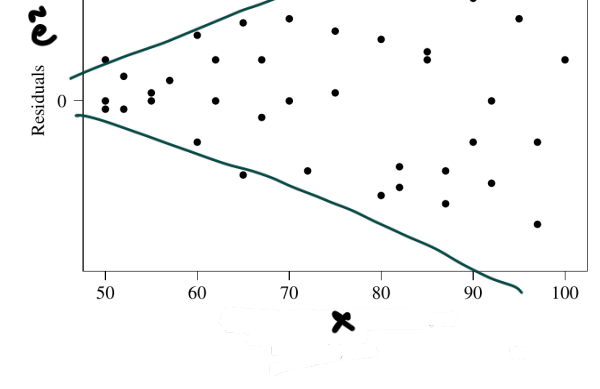
ok!
• no patterns
• constant dispersion



No! systematic behaviors
quadratic trend is suggesting that we should include x^2 in the model



other examples:
 y depends linearly on a covariate not included in the model



No! heteroscedastic: the variance increases with x



other examples

(2) NORMALITY ASSUMPTION

we can use the studentized residuals $R_i \sim N(0, 1)$

- histogram of R_i vs. normal density (but it is not so simple to identify deviations)
- empirical cumulative distribution function (ECDF) vs CDF Φ of a $N(0, 1)$
- normal Q-Q plot (quantile-quantile plot)

Recall: ECDF. If I observe a sample u_1, \dots, u_n from U with $F_U(t) = P(U \leq t)$

the empirical CDF is $\hat{F}_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}(u_i \leq t) = \frac{\text{number of observations } \leq t}{n}$

$\hat{F}_n(t)$ is an unbiased estimator of $F_U(t)$.

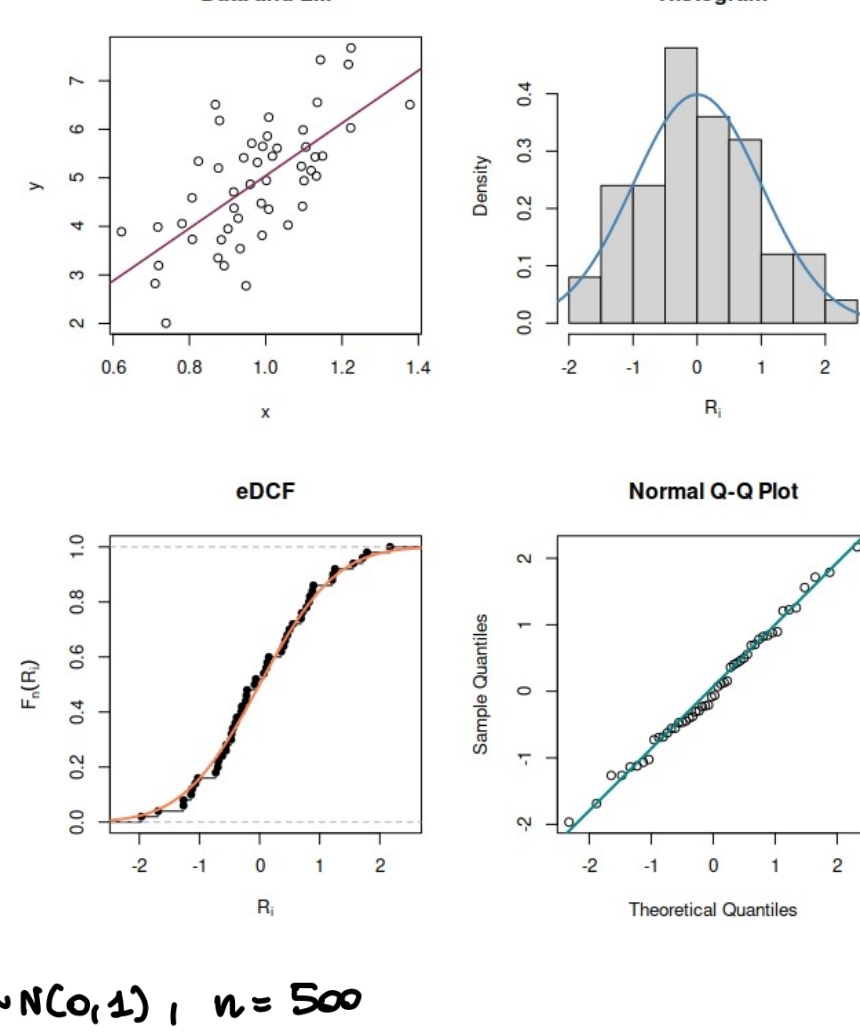
Recall: normal Q-Q plot. It is a plot of the empirical quantiles vs. the theoretical quantiles of a $N(0, 1)$. If the observed sample comes from a $N(0, 1)$ distribution, the points should lie on the straight line $y=x$.

(empirical quantiles = ordered sample $u_{(1)} \leq u_{(2)} \leq \dots \leq u_{(n)}$)
theoretical quantiles = Φ^{-1} inverse of the CDF of a $N(0, 1)$)

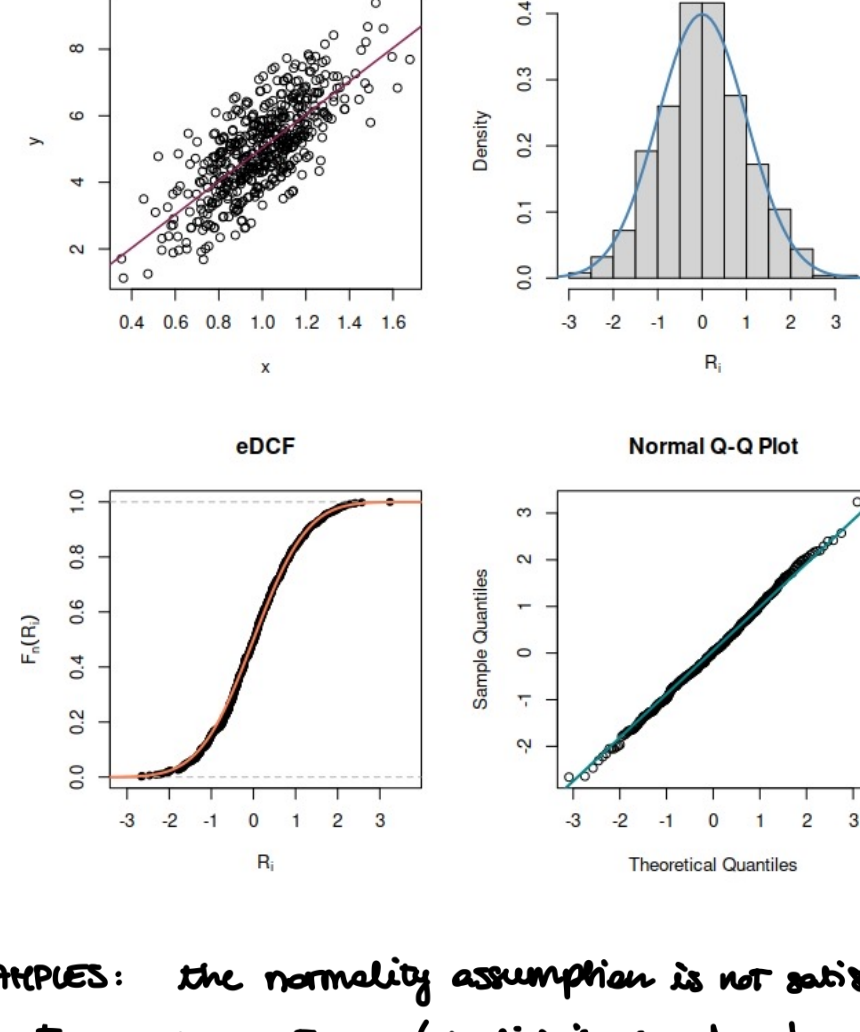
EXAMPLES: $y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \quad i=1, \dots, n$

the normality assumption is satisfied $\Rightarrow R_i \sim N(0, 1)$

$\epsilon_i \sim N(0, \sigma^2)$ and $n = 500$

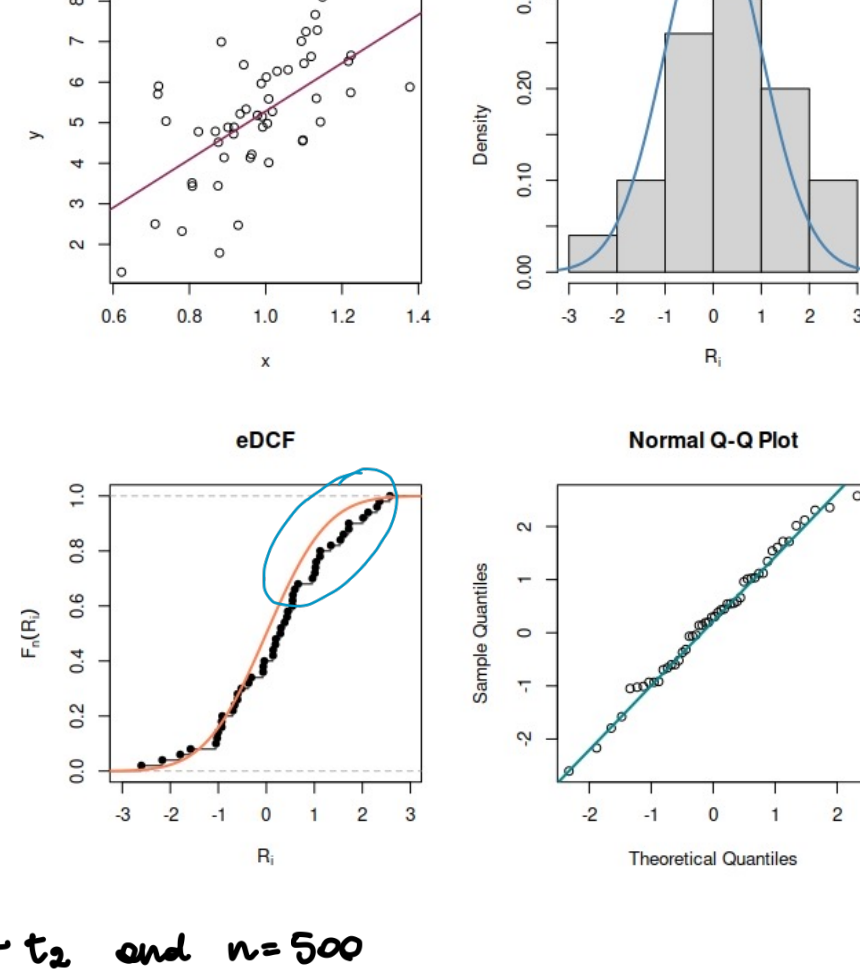


$\epsilon_i \sim N(0, 1), \quad n = 500$



EXAMPLES: the normality assumption is not satisfied

$\epsilon_i \sim t_2$ and $n = 50$ (t distribution has heavier tails)



$\epsilon_i \sim t_2$ and $n = 500$

