

Regresní model procházející počátkem

Existují případy, kdy přípustný model vyžaduje $\beta_0 = 0$, tj.

$$Y_i = \beta_1 x_i + e_i, \quad i = 1, \dots, n,$$

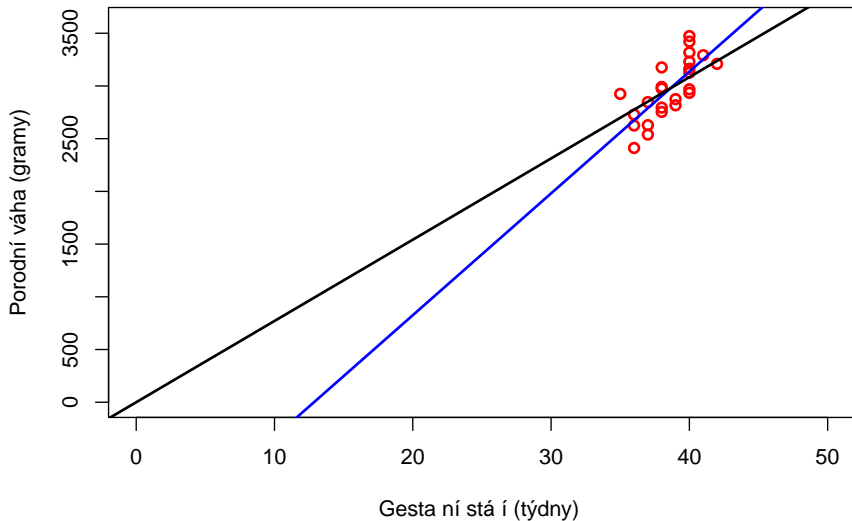
např.

- je to předem známo na základě fyzikálních úvah ($EY_0 = \beta_0 = 0$), potom nemá smysl odhadovat β_0 , obecně to snižuje přesnost odhadu σ^2 a tím i β_1
- na začátku předpokládáme $\beta_0 \neq 0$ a t-test nezamítne $H_0 : \beta_0 = 0$, potom β_0 může být z modelu odstraněno

POZNÁMKA 2.15

- v praxi často není jisté, že model platí i blízko počátku
- část statistiků trvá na přítomnosti interceptu v modelu, i když je nevýznamný
- položit $\beta_0 = 0$ apriorně může být chybné, i když $EY_0 = 0$, pokud totiž nevíme jistě, že model je lineární na okolí 0, volba $\beta_0 = 0$ může vést k vychýleným odhadům β_1 , pokud jsou nezávislé proměnné daleko od 0

Příklad: porodní váha a gestační stáří



Odhady a testy v případě $\beta_0 = 0$

- LSE parametru β_1 dostaneme minimalizací $S = \sum_{i=1}^n (y_i - \beta_1 x_i)^2$ ve tvaru $\hat{\beta}_1 = \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2}$
- pokud e_1, \dots, e_n i.i.d. $N(0, \sigma^2)$, potom

$$E\hat{\beta}_1 = \beta_1, \quad \text{Var}\hat{\beta}_1 = \frac{\sigma^2}{\sum_{i=1}^n x_i^2}, \quad \hat{\beta}_1 \sim N\left(\beta_1, \frac{\sigma^2}{\sum_{i=1}^n x_i^2}\right),$$

$$s_n^2 = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n-1} = \frac{SSE}{n-1} \text{ je nestranný odhad } \sigma^2, \quad \frac{SSE}{\sigma^2} \sim \chi^2(n-1) \text{ a nezávisí na } \hat{\beta}_1$$

- $H_0 : \beta_1 = 0$ lze otestovat pomocí $T = \frac{\hat{\beta}_1}{\frac{s_n}{\sqrt{\sum_{i=1}^n x_i^2}}} \sim t(n-1)$

- 100(1 - α)% IS pro β_1 je $\left(\hat{\beta}_1 \pm t_{1-\alpha/2}(n-1) s_n / \sqrt{\sum_{i=1}^n x_i^2}\right)$

POZNÁMKA 2.16

- zatím vše podobné jako pro případ $\beta_1 \neq 0$
- rozdíl ale bude v tabulce ANOVA a R^2 statistice, neplatí totiž rozklad $SST = SSR + SSE$
- odvodíme nový rozklad (platí v obou modelech, dokážeme jen pro $\beta_0 = 0$)

VĚTA 2.6

V modelu s $\beta_0 = 0$ platí $\sum_{i=1}^n y_i^2 = \sum_{i=1}^n \hat{y}_i^2 + \sum_{i=1}^n (y_i - \hat{y}_i)^2$.

Důkaz.

Pokud vezmeme $\sum_{i=1}^n y_i^2$ jako míru variability v datech, analogie R^2 bude: $R^2 = \frac{\sum_{i=1}^n \hat{y}_i^2}{\sum_{i=1}^n y_i^2}$

potom

$$1 - R^2 = \frac{\sum_{i=1}^n y_i^2 - \sum_{i=1}^n \hat{y}_i^2}{\sum_{i=1}^n y_i^2} = \frac{\sum_{i=1}^n \hat{e}_i^2}{\sum_{i=1}^n y_i^2}, \quad \text{a definicí} \quad F = \frac{(n-1)R^2}{1 - R^2}$$

dostaneme

$$F = \frac{\sum_{i=1}^n \hat{y}_i^2}{\frac{1}{n-1} \sum_{i=1}^n (y_i - \hat{y}_i)^2} = \frac{\hat{\beta}_1^2 \sum_{i=1}^n x_i^2}{s_n^2} = T^2$$

vztah mezi R^2, F, T je tedy stejný jako pro $\beta_0 \neq 0$

POZNÁMKA 2.17

Tato definice R^2 se ale moc nepoužívá, protože neumožňuje přímé srovnání modelů *bez* a *s* interceptem

$$\beta_0 = 0: R^2 = 1 - \frac{SSE}{\sum_{i=1}^n y_i^2}$$

$$\beta_0 \neq 0: R^2 = 1 - \frac{SSE}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

obecně ale $\sum_{i=1}^n (y_i - \bar{y})^2 < \sum_{i=1}^n y_i^2$, R^2 modelu s $\beta_0 = 0$ tak bude větší než R^2 modelu s $\beta_0 \neq 0$ (i když jsou jejich SSE srovnatelné)

- definice vhodné R^2 pro $\beta_0 = 0$ vyvolává jistou kontroverzi a ex. několik verzí
- možná volba je

$$R^2 = \varrho^2(\mathbf{y}, \hat{\mathbf{y}}),$$

kde $\hat{\mathbf{y}} = (\hat{y}_1, \dots, \hat{y}_n)$ (platí i pro $\beta_0 \neq 0$)

- další možnost: srovnat modely pomocí s_n^2

Tabulka ANOVA pro $\beta_0 = 0$.

Source	df	SS	MS	F
Regression	1	$SSR = \sum_{i=1}^n \hat{y}_i^2$	$MSR = \frac{SSR}{1}$	$\frac{SSR}{s_n^2}$
Residual	$n - 1$	$SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2$	$MSE = \frac{SSE}{n-1} = s_n^2$	
Total	n	$SST = \sum_{i=1}^n y_i^2$		

$$R^2 = \varrho^2(\mathbf{y}, \hat{\mathbf{y}})$$

PŘÍKLAD 2.4 (Porodní váha a gestační stáří)

```
mod <- lm(Weight ~ Age)
summary(mod)
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -1485.0      852.6   -1.742   0.0955 .
## Age           115.5       22.1    5.228 3.04e-05 ***
## Residual standard error: 192.6 on 22 degrees of freedom
## Multiple R-squared:  0.554,    Adjusted R-squared:  0.5338
## F-statistic: 27.33 on 1 and 22 DF,  p-value: 3.04e-05
```

```
y.mod <- predict(mod)
cor(y.mod,Weight)^2
## 0.5540268
```

Model bez interceptu:

```
mod.bez <- lm(Weight ~ Age - 1)
summary(mod.bez)

## Coefficients:
##      Estimate Std. Error t value Pr(>|t|)
## Age    77.081      1.063   72.51  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

## Residual standard error: 200.9 on 23 degrees of freedom
## Multiple R-squared:  0.9956,      Adjusted R-squared:  0.9955
## F-statistic: 5258 on 1 and 23 DF,  p-value: < 2.2e-16

anova(mod.bez)

## Analysis of Variance Table
##
## Response: Weight
##      Df    Sum Sq  Mean Sq F value    Pr(>F)
## Age      1 212270368 212270368  5257.6 < 2.2e-16 ***
## Residuals 23   928596    40374
##

y.mod.bez <- predict(mod.bez)

cor(y.mod.bez,Weight)^2

## 0.5540268
```


Predikce

Jakmile máme model, často bývá cílem odhadnout hodnoty veličiny Y_0 pro nové x_0 , které není v původních datech.

Budeme uvažovat dva typy predikce:

- 1) predikce střední hodnoty $\mu_0 = E[Y_0]$ v bodě x_0 ,
- 2) predikce hodnoty nového pozorování Y_0 v bodě x_0 .

Pro oba typy použijeme bodový odhad

$$\hat{Y}_0 = \hat{\beta}_0 + \hat{\beta}_1 x_0,$$

intervalové odhady se ale budou lišit.

Ad 1) $\mu_0 = \beta_0 + \beta_1 x_0$ je vlastně parametr, lze pro něj odvodit IS (za předpokladu normality chyb).

Spočteme $\text{Var}(\hat{Y}_0)$:

Shrnutí: $100(1 - \alpha)\%$ IS pro μ_0

$$\hat{Y}_0 \pm t_{1-\frac{\alpha}{2}}(n-2) \cdot \hat{\sigma}(\hat{Y}_0), \quad \text{kde} \quad \hat{\sigma}^2(\hat{Y}_0) = s_n^2 \left[\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right]$$

$\hat{\sigma}(\hat{Y}_0)$ se obvykle nazývá **standardní chyba predikce v bodě x_0**

POZNÁMKA 2.18

Z tvaru IS je vidět, že bude nejkratší pro $x_0 = \bar{x}$ a s rostoucí vzdáleností $|x_0 - \bar{x}|$ se prodlužuje.

- Speciálně potom čím dále jsme od oblasti, kde jsou naše data x , tím méně spolehlivé jsou naše predikce.
- Je třeba opatrnosti při predikci hodnot Y mimo interval $(\min x_i, \max x_i)$.

Ad 2)

Intervalové odhady pro Y_0 nejsou IS, protože Y_0 není parametr \longleftrightarrow intervaly predikce.

Potřebujeme rozptyl $Y_0 - \hat{Y}_0$:

Shrnutí: $100(1 - \alpha)\%$ interval predikce pro Y_0

$$\hat{Y}_0 \pm t_{1-\frac{\alpha}{2}}(n-2) \cdot s_p, \quad \text{kde} \quad s_p = s_n \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}}}$$

POZNÁMKA 2.19

Přesnost predikce

- a) roste s rostoucím n a rostoucím rozsahem x naměřeným pomocí S_{xx} ,
- b) klesá s rostoucím $|x_0 - \bar{x}|$.

Pokud je možno předem zvolit x_1, \dots, x_n , lze přesnost predikce zvýšit volbou dostatečně rozptýlených hodnot x .

To ale může zvyšovat R^2 a někdy vést k horšímu modelu.

⇒ **základní rozpor v regresní analýze:**

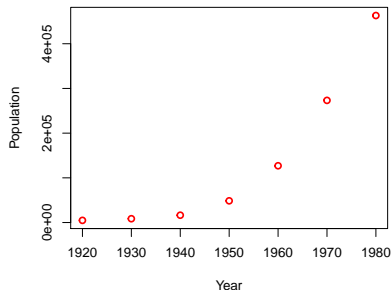
- dobrý model nemusí poskytovat dobré predikce,
- dobré predikce mohou vycházet z méně přesných modelů.

POZNÁMKA 2.20

- odvozené výsledky platí za předpokladu normality chyb
- za podmínek regularity jsou ale odhady $\hat{\beta}_0, \hat{\beta}_1$ asymptoticky normální, tzn. **IS pro EY_0** budou použitelné pro velká n i pro nenormální chyby
- **IP pro Y_0** ale závisí na normalitě chyb i pro velká n , tzn. mohou tedy být nepřesné pro nenormální chyby

PŘÍKLAD 2.5 (Clark County population data)

x	Year	Population
0	1920	4859
1	1930	8539
2	1940	16414
3	1950	48589
4	1960	127016
5	1970	273288
6	1980	463087



1. lineární model pro původní data: $Y_i = \beta_0 + \beta_1 x_i + e_i$

$$\hat{\beta}_0 = -81\,328, \quad p_{val} = 0.22, \quad \hat{\beta}_1 = 71\,957, \quad p_{val} = 0.007, \quad R^2 = 0.80, \quad F = 19.96$$

x	Year	Population	Fitted.value
0	1920	4859	-81328
1	1930	8539	-9371
2	1940	16414	62585
3	1950	48589	134542
4	1960	127016	206498
5	1970	273288	278455
6	1980	463087	350411

Predikce pro rok 1990:

$$\hat{y}_{1990} = \hat{\beta}_0 + \hat{\beta}_1 \cdot 7 = 422\,368$$

a) 95% IS: (237 233, 607 502)

b) 95% IP: (135 559, 709 177)

Skutečná hodnota v roce 1990: 768 203

2. lineární model pro log-transformovaná data: $\log(Y_i) = \beta_0 + \beta_1 x_i + e_i$

$$\hat{\beta}_0 = 8.33, \quad p_{val} < 10^{-4}, \quad \hat{\beta}_1 = 0.809, \quad p_{val} < 10^{-4}, \quad R^2 = 0.991, \quad F = 550.9$$

Predikce pro rok 1990 na log. škále:

$$\widehat{\log y}_{1990} = \hat{\beta}_0 + \hat{\beta}_1 \cdot 7 = 14.0004$$

a) 95% IS: (13.604, 14.397)

b) 95% IP: (13.387, 14.614)

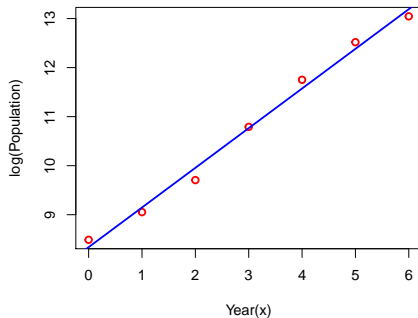
Predikce pro rok 1990 na původní škále:

$$\hat{y}_{1990} = 1\,203\,161$$

a) 95% IS: (809 576, 1 788 092)

b) 95% IP: (651 269, 2 222 733)

Skutečná hodnota v roce 1990: 768 203



Intervaly predikce v R

```
mod.lin <- lm(pop ~ year)

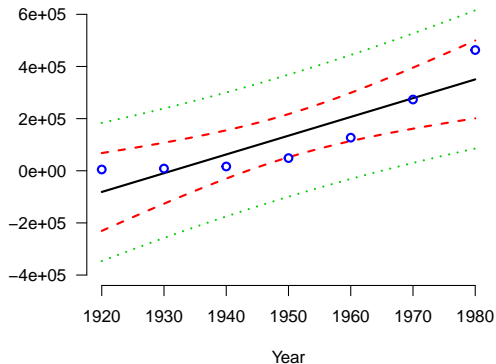
new <- data.frame(year=seq(0,6,0.1))

# predikce y
y.hat <- predict(lm(pop ~ year), new)

# 95% intervaly spolehlivosti
CI<-predict(mod.lin, new, interval = "confidence")

# 95% intervaly predikce
PI<-predict(mod.lin, new, interval = "prediction")

# obrazek
matplot(new$year, cbind(CI, PI[, -1]),
        lty = c(1,2,2,3,3), col = c(1,2,2,3,3),
        type = "l", lwd=2, ylab="", xlab = "Year",
        axes = FALSE, ylim=c(-400000,600000))
axis(side=1, at=0:6, labels=Year)
axis(side=2, las=2)
points(pop ~ year, col="blue", lwd=2)
```



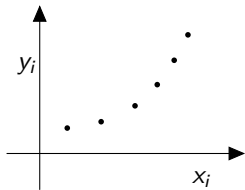
Ověření adekvátnosti modelu

- důležitá součást analýzy
- mělo by předcházet interpretaci modelu případně přijímání závěrů založených na modelu
- výsledky odvozeny za předpokladu **linearity modelu**, případně **normality chyb**

Základní procedury:

1) Prozkoumání scatter plotu dvojic (x_i, y_i)

Např.



může indikovat, že lepší model bude

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + e_i.$$

(může být zavádějící)

2) Analýza hodnot testovacích statistik

- např. malá hodnota R^2 společně s významnou hodnotou t -statistiky pro parametr β_1 naznačuje, že skutečný model obsahuje i jiné proměnné x
- velká hodnota R^2 a významná t -statistika ale samo o sobě neznamená, že je model lineární.

3) Obrázky reziduí

- efektivní diagnostický nástroj
- rezidua odhadují, kolik variability v datech zůstane po odstranění lineární části v x
- dá se očekávat, že budou užitečné pro detekci odchylek od normality

Ad 3) - Analýza reziduí

- intuitivně, pokud je náš model správný, měla by se rezidua chovat jako náhodný výběr z $N(0, \sigma^2)$
- pokud se tak nechovají, může to znamenat neadekvátnost modelu
- ukážeme grafické nástroje, začneme ale vlastnostmi reziduí

VĚTA 2.7

Nechť \hat{e}_i jsou rezidua modelu (*) odhadnutého metodou nejmenších čtverců. Potom platí:

- 1) $E(\hat{e}_i) = 0, \quad i = 1, \dots, n$
- 2) $\text{Var}(\hat{e}_i) = \sigma_{\hat{e}_i}^2 = \sigma^2 \left[1 - \left(\frac{1}{n} + \frac{(x_i - \bar{x})^2}{S_{xx}} \right) \right] \quad (\approx \sigma^2 \text{ pro velká } n)$
- 3) $\text{Cov}(\hat{e}_i, \hat{e}_j) = -\sigma^2 \left[\frac{1}{n} + \frac{(\bar{x} - x_i)(\bar{x} - x_j)}{S_{xx}} \right]$
- 4) $\text{Cov}(\hat{e}_i, \hat{Y}_i) = 0, \quad i = 1, \dots, n$
- 5) Pokud jsou e_1, \dots, e_n i.i.d. $N(0, \sigma^2)$, potom platí: $\hat{Z}_i = \frac{\hat{e}_i}{\sigma_{\hat{e}_i}} \sim N(0, 1).$

Důkaz.

POZNÁMKA 2.21

- bod 3) věty $\Rightarrow \text{Cov}(\hat{e}_i, \hat{e}_j) \approx 0$ pro velké n
- pokud jsou tedy e_i i.i.d. $N(0, \sigma^2)$, měla by se *standardizovaná rezidua* $\hat{Z}_i = \frac{\hat{e}_i}{\hat{\sigma}_{e_i}}$ chovat pro velké n jako náhodný výběr z $N(0, 1)$ rozdělení
- budeme potřebovat odhad σ^2 pro výpočet \hat{Z}_i
- nejznámější procedura: odhadnout σ^2 pomocí s_n^2 , potom

$$\hat{r}_i = \frac{\hat{e}_i}{s_n \sqrt{1 - \left(\frac{1}{n} + \frac{(x_i - \bar{x})^2}{S_{xx}} \right)}} \quad \text{studentizovaná rezidua}$$

- pro velká n by se opět měla \hat{r}_i chovat jako náh. výběr z $N(0, 1)$.

POZNÁMKA 2.22

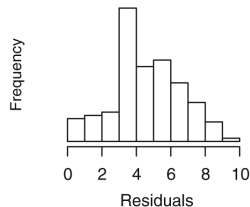
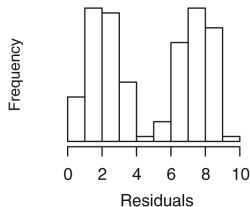
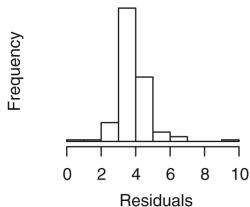
- \hat{e}_i, \hat{r}_i se užívají pro grafickou analýzu
- jiná třída reziduí – **PRESS rezidua** (negrafické metody zkoumání reziduí):
ozn. $\hat{\beta}_{0(-i)}, \hat{\beta}_{1(-i)}$ odhady parametrů β_0, β_1 , pokud je vynecháno i -té pozorování
pak i -té PRESS reziduum je definováno jako

$$\hat{e}_{(-i)} = y_i - \hat{y}_{(-i)}, \quad \text{kde } \hat{y}_{(-i)} = \hat{\beta}_{0(-i)} + x_i \hat{\beta}_{1(-i)}.$$

(podrobněji se jim budeme věnovat později)

Grafy reziduí

1) histogram reziduí

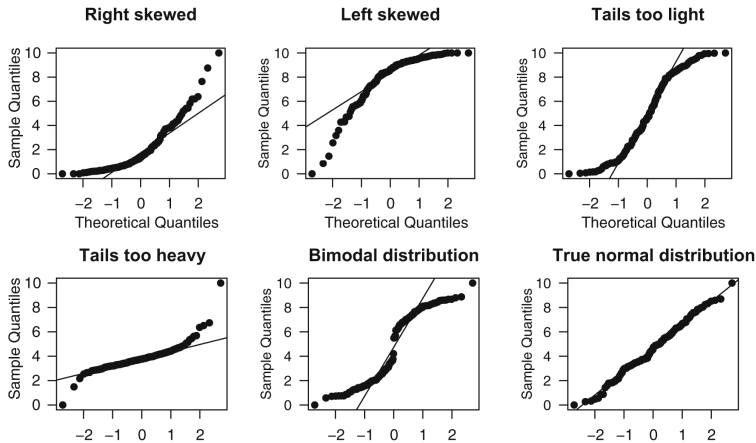


2) kvantilový graf (Q-Q plot) studentizovaných reziduí

- seřadíme dle velikosti:

$$\hat{r}_{(1)} \leq \hat{r}_{(2)} \leq \dots \leq \hat{r}_{(n)} \quad \text{a vyneseme oproti} \quad \Phi^{-1}\left(\left(i - \frac{1}{2}\right)\frac{1}{n}\right), i = 1, \dots, n$$

- body by měly ležet přibližně na přímce ($E(r_{(i)}) \approx \Phi^{-1}\left(\left(i - \frac{1}{2}\right)\frac{1}{n}\right)$ pro normální chyby)
- použití:** ověření normality, detekce odlehlých pozorování

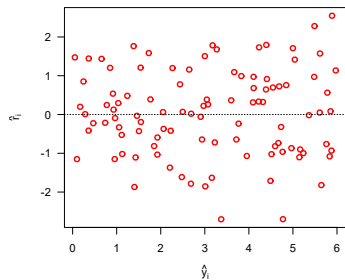


3) studentizovaná rezidua vs. jednotlivé vysvětlující proměnné x

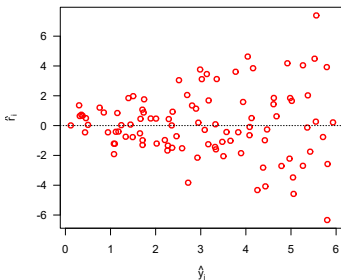
\hat{r}_i nezávisí na σ , graf $\hat{r}_i \times x_i$ lze použít pro detekci nelinearity nebo nekonstantního rozptylu

4) studentizovaná rezidua \hat{r}_i vs. predikované hodnoty \hat{y}_i

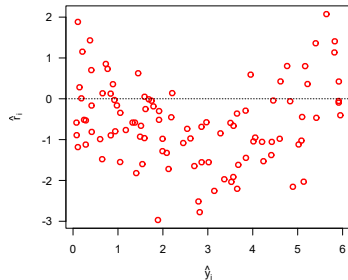
- $\text{Cov}(\hat{e}_i, \hat{Y}_i) = 0$, tedy \hat{r}_i a \hat{Y}_i by měly být nekorelované, pokud platí model (*)
- tzn. graf $\hat{r}_i \times \hat{y}_i$ by měl být náhodně rozptýlený kolem osy x
- navíc \hat{r}_i by měla ležet v $(-3, 3)$ ($\hat{r}_i \approx N(0, 1)$)



(a)



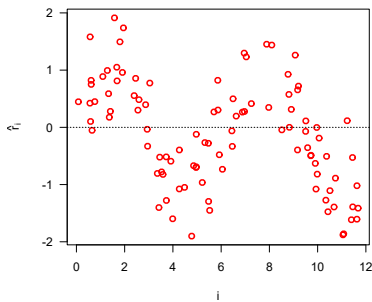
(b)



(c)

5) studentizovaná rezidua vs. pořadí pozorování

možná detekce řadové korelace mezi pozorováními



Obrázek: Studentizovaná rezidua \hat{r}_i proti pořadí pozorování i .