

Kazalo

1	NELINEARNOST	1
1.1	Polinomska regresija	1
1.2	Regresija zlepkov	11
1.2.1	Bazne funkcije	12
1.2.2	Linearni zleпки	12
1.2.3	Kubični zleпки	14
1.2.4	Naravni zleпки	14
1.2.5	Primer: KORUZA (nadaljevanje)	16
2	VAJE	25
2.1	Telesna masa in višina žensk	25
2.2	Plača	25

1 NELINEARNOST

Linearnost odvisnosti odzivne spremenljivke od napovedne spremenljivke ob upoštevanju ostalih spremenljivk v modelu se v praksi velikokrat pokaže kot precej slaba aproksimacija dejanske odvisnosti. Obstaja več načinov modeliranja nelinearnosti v kontekstu linearnih modelov. Najpreprostejši sta polinomska regresija in regresija po odsekih (*step function regression* in *piecewise regression*), kompleksnejše metode so regresija zlepkov (*regression splines*), glajeni zleпки (*smoothing splines*), lokalna regresija (*local regression*) in posplošeni aditivni modeli (*Generalized Additive Models*, GAM). V tem poglavju bomo predstavili polinomsko regresijo in regresijo zlepkov.

1.1 Polinomska regresija

Zgodovinsko gledano predstavlja polinomska regresija najstarejši način modeliranja nelinearne odvisnosti odzivne od napovedne spremenljivke. Osnovni linearni model

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad (1)$$

v tem primeru razširimo v polinomom stopnje p , $p \geq 2$:

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \dots + \beta_p x_i^p + \varepsilon_i. \quad (2)$$

Ob upoštevanju $x = x_1, x^2 = x_2, \dots, x^p = x_p$, lahko izraz (2) zapišemo kot model, ki vključuje več številskih napovednih spremenljivk:

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}.$$

Spremenljivke x_1, \dots, x_p so med seboj odvisne (multikolinearnost), kar pa ne predstavlja večjih težav, saj običajno na podlagi takega modela ne preverjamo ničelnih domnev o posameznih

parametrih modela, bolj nas zanimajo napovedi. Parametri $\beta_0, \beta_1, \dots, \beta_p$ so v linearnem odnosu z y , nimajo pa vsebinskega pomena.

Če želimo preveriti, ali je linearna odvisnost y od x upravičena, preverjamo sestavljeno ničelno domnevo $H_0 : \beta_2 = \beta_3 = \dots = \beta_p = 0$. To naredimo z F-testom za gnezdene modele.

Pri modeliranju s polinomske regresijo se v praksi skušamo omejiti na polinome nižjih stopenj, $p = 2$ do 4. Pri polinomih stopnje več kot 4 hitro pride do preprileganja podatkov, še posebej na robovih prostora napovedne spremenljivke. Namesto polinomov višjih stopenj je v določenih primerih bolje uporabiti nelinearne regresijske modele, pri katerih se parametri dajo vsebinsko interpretirati ali pa regresijo zlepkov.

Primer: koruza

Za rezultate bločnega poskusa s koruzo v letu 1990 (KORUZA.txt) analizirajmo, kako je pridelek koruze (kg/ha) odvisen od gostote setve. Zanima nas optimalen pridelek koruze, optimalna gostota setve in njuna 95 % intervala zaupanja. Poskus je bil zasnovan kot bločni poskus v 3 ponovitvah (blokih), v poskusu je bilo 15 različnih gostot setve. Znotraj vsakega bloka (dela njive) so bile enkrat ponovljene vse gostote setve. Ker tudi blok lahko vpliva na pridelek (npr. zaradi različnih rastnih pogojev), bomo vpliv gostote setve na pridelek modelirali ob upoštevanju vpliva bloka. V tem primeru blok vključimo v model kot opisno spremenljivko, njen vpliv pa je slučajen (o tem več v poglavju o mešanih linearnih modelih).

```
> koruza<-read.table(file="KORUZA.txt", header = TRUE)
> str(koruza)

'data.frame':      45 obs. of  5 variables:
 $ blok      : int  1 1 1 1 1 1 1 1 1 1 ...
 $ gostsetve : num  65.2 23.6 129.9 50.5 47.6 ...
 $ gostvznika: num  51.9 23.5 123.8 49.5 41.3 ...
 $ prid.ha   : int  5880 2936 6962 5152 4129 720 5219 6481 3508 6719 ...
 $ prid.rast : num  0.09 0.124 0.054 0.102 0.087 0.145 0.073 0.038 0.101 0.071 ...

> summary(koruza[,c("gostsetve", "prid.ha")])

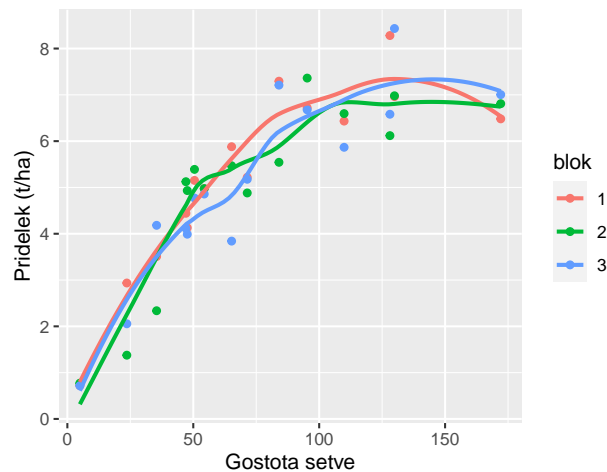
      gostsetve      prid.ha
Min.   : 4.967   Min.     : 717
1st Qu.: 47.080   1st Qu.: 4129
Median : 65.230   Median : 5176
Mean    : 74.632   Mean     : 5095
3rd Qu.: 109.900   3rd Qu.: 6595
Max.    : 172.100   Max.     : 8433
```

Spremenljivka `prid.ha` je izražena v kg/ha, zaradi lažje interpretacije jo pretvorimo v t/ha, spremenljivko `blok` pa spremenimo v `factor`:

```
> koruza$prid1.ha <- koruza$prid.ha/1000
> koruza$blok <- factor(koruza$blok)
```

Slika 1 prikazuje odvisnost pridelka koruze od gostote setve in od bloka s pripadajočimi gladilniki.

```
> library(ggplot2)
> ggplot(data=koruza, aes(x=gostsetve, y=prid1.ha, col=blok)) +
+   geom_point() + geom_smooth(se=FALSE) +
+   ylab("Pridelek (t/ha)") +
+   xlab("Gostota setve")
```



Slika 1: Odvisnost pridelka (t/ha) od gostote setve in od bloka z dodanim gladilnikom

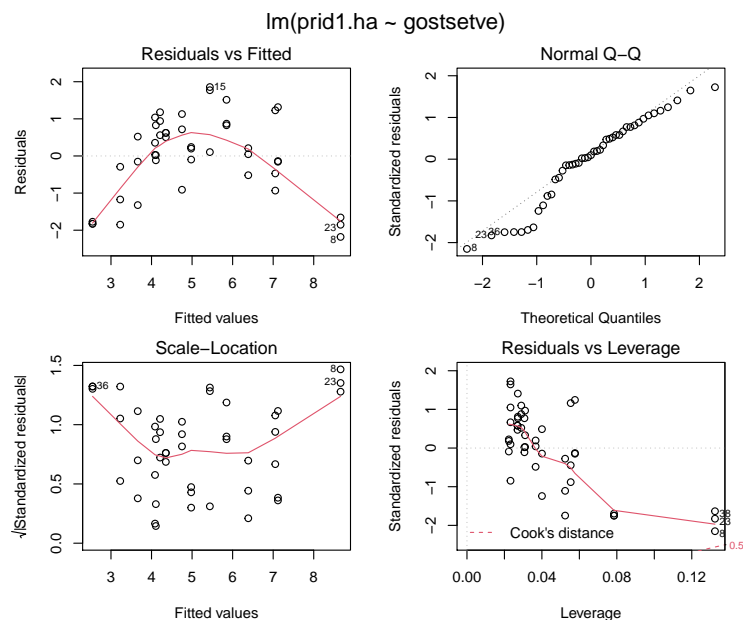
Slika 1 ter vsebinski premislek nakazujejo, da pridelek ni linearno odvisen od gostote setve. Vidimo, da je odvisnost za vse tri bloke zelo podobna. Naredimo najprej neustrezeni linearni model in hkrati pogledjmo, ali lahko vpliv blokov zanemarimo.

```
> model.lin <- lm(prid1.ha ~ gostsetve, data=koruza)
> model.lin.blok <- lm(prid1.ha ~ blok + gostsetve, data=koruza)
> anova(model.lin, model.lin.blok)
```

Analysis of Variance Table

```
Model 1: prid1.ha ~ gostsetve
Model 2: prid1.ha ~ blok + gostsetve
  Res.Df  RSS Df Sum of Sq    F Pr(>F)
1     43 50.941
2     41 50.187  2   0.75458 0.3082 0.7364
```

Modela `model.lin` in `model.lin.blok` sta ekvivalentna, zato nadaljujemo s preprastejšim modelom, ki kot napovedno spremenljivko vključuje samo `gostsetve`.



Slika 2: Ostanki za `model.lin`

Ostanki kažejo, da `model.lin` ni sprejemljiv, na Grafu 1 je gladilnik v obliki parabole, zato model dopolnimo s kvadratnim členom:

```
> model.kvad <- lm(prid1.ha ~ gostsetve + I(gostsetve^2), data=koruza)
> # enak rezultat dobimo z uporabo funkcije poly
> # model.kvad.1 <- lm(prid1.ha ~ poly(gostsetve, degree=2, raw=TRUE), data=koruza)
```

V zapisu `lm` modela uporabimo funkcijo `I()`, ki zagotovi, da izraz `gostsetve^2` določa regresor v linearnem modelu, znak za potenco, kot tudi ostali aritmetični operatorji (`*`, `/`, `+`, `-`), ima v formuli modela poseben pomen in s tem je ta pomen omejen na potenciranje `gostsetve`.

V `model.kvad` sta regresorja korelirana in posledično je VIF vrednost visoka, kar v tem primeru ignoriramo.

```
> library(car)
> vif(model.kvad)

gostsetve I(gostsetve^2)
13.16808      13.16808
```

Naredimo primerjavo `model.lin` in `model.kvad`. Preverjamo ničelno domnevo, da sta modela ekvivalentna. F -test za dva gnezdena modela izvedemo s funkcijo `anova`

```
> anova(model.kvad)
```

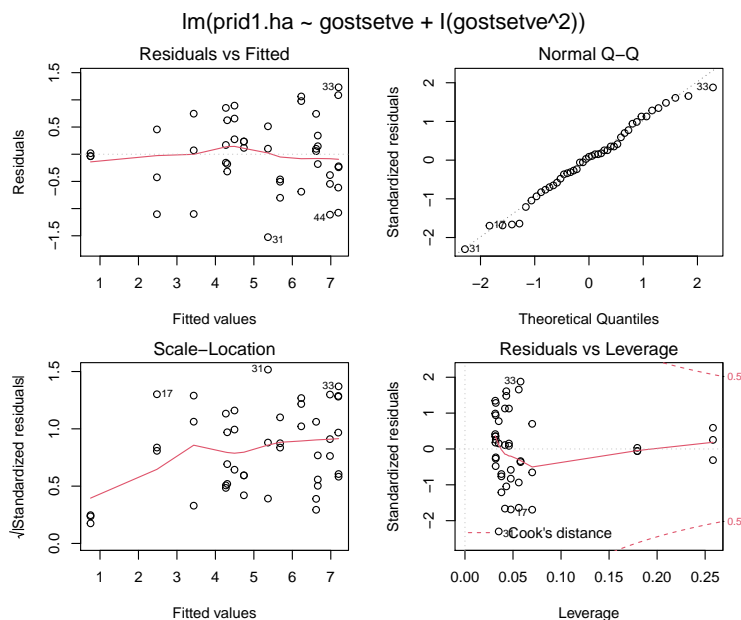
Analysis of Variance Table

Response: prid1.ha

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
gostsetve	1	115.495	115.495	253.488	< 2.2e-16 ***
I(gostsetve^2)	1	31.805	31.805	69.807	1.798e-10 ***
Residuals	42	19.136	0.456		

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Model z dodanim kvadratnim členom je statistično značilno boljši od modela brez kvadratnega člena ($F = 69.8, p < 0.0001$).



Slika 3: Ostanke za model.kvad

Naraščajoč gladilnik na Sliki 3 levo spodaj je pretežno posledica treh podatkov pri zelo nizki vrednosti napovedanega pridelka, če te točke odmislimo, težav z nekonstantno varianco ni videti in lahko rečemo, da je `model.kvad` sprejemljiv. Poleg tega nas nekonstantna varianca ne skrbi, ker ni cilj modeliranja testiranje ničelne domneve o statističnem vplivu gostote na pridelok, temveč iskanje optimalne gostote setve. Videli smo, da na nepristranskost ocen parametrov nekonstantna varianca ne vpliva (pri dokazovanju nepristranskosti ocen parametrov modela nismo uporabili predpostavke o konstantni varianci).

```
> coefficients(model.kvad)

(Intercept)      gostsetve I(gostsetve^2)
0.2491926865  0.1035883074 -0.0003853948

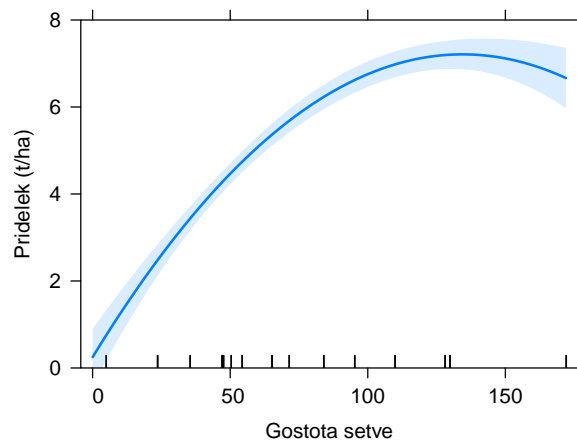
> summary(model.kvad)$r.squared

[1] 0.8850243
```

Napišimo enačbo parabole:

$$\hat{y} = 0.24919 + 0.10359x + (-0.00039)x^2.$$

```
> library(effects)
> plot(Effect(c("gostsetve"), model.kvad, xlevels=list(gostsetve=seq(0, 172, 2))),
+      ci.style="bands", xlab="Gostota setve", ylab="Pridelek (t/ha)",
+      main="", ylim=c(0,8))
```



Slika 4: Odvisnost pridelka (t/ha) od gostote setve in parabola izračunana z `model.kvad` ter 95 % intervali zaupanja za povprečne napovedi

Optimalna gostota setve

Izračunajmo optimumalno gostoto setve, optimalni pridelek in pripadajoči 95 % IZ. Optimum izračunamo z odvajanjem kvadratne enačbe po x in enačenjem odvoda z 0:

$$b_1 + 2b_2x = 0, \quad x = -\frac{b_1}{2b_2}$$

```
> opt<--0.5*coefficients(model.kvad)[2]/coefficients(model.kvad)[3]
> cat("Optimalna gostota =", opt)
```

Optimalna gostota = 134.3925

```
> # Napoved in interval zaupanja za pridelek pri optimalni gostoti
> gostsetve.x<-data.frame(gostsetve=opt)
> povp.napoved.pridelek<-predict(model.kvad,gostsetve.x, interval="confidence")
> round(data.frame(cbind(gostsetve.x,povp.napoved.pridelek)), 2)
```

	gostsetve	fit	lwr	upr
gostsetve	134.39	7.21	6.87	7.55

Interpretacija rezultatov: pri optimalni gostoti 134.39 je pričakovana vrednost pridelka 7.21 t/ha, pripadajoči 95 % IZ pa je (6.87 t/ha, 7.55 t/ha).

Intervalna ocena za optimalno gostoto setve

Optimum je izračunan kot razmerje dveh normalno porazdeljenih slučajnih spremenljivk in je tudi slučajna spremenljivka. Zanima nas njen interval zaupanja, za to rabimo pripadajočo varianco. Tega ne znamo dobiti analitično, uporabimo lahko eno izmed metod samovzorčenja.

Z **neparametričnim bootstrap pristopom** (Efron, 1979) iz osnovnega vzorca velikosti n tvorimo t. i. **bootstrap vzorce**. Vsak bootstrap vzorec ima n enot. Enote (s pripadajočimi vrednostmi odzivne in napovednih spremenljivk) vzorčimo z enostavnim slučajnim vzorčenjem s ponavljanjem. Takemu načinu samovzorčenja v kontekstu linearnih modelov pravimo **samovzorčenje primerov** (*case resampling*). Tvorimo R bootstrap vzorcev, R je veliko število (1000 in več). Za vsak bootstrap vzorec izračunamo vzorčno oceno iskane statistike. Na osnovi R bootstrap vzorcev dobimo njeno **bootstrap vzorčno porazdelitev**. 95 % **centilni bootstrap interval zaupanja** je določen z 2.5 in 97.5 centilom te porazdelitve.

```
> R <- 1000
> b <- data.frame()
> for (i in c(1:R)){
+   vzorec <- sample(c(1:dim(koruza)[1]),replace=TRUE)
+   koruza1<-koruza[vzorec,]
```

```
+ b <- rbind(b, coefficients(lm(prid1.ha ~ gostsetve + I(gostsetve^2), data=koruza1))
+ }
> names(b) <-c("b0", "b1", "b2")
> b$opt <- -b$b1/(2*b$b2)
> mean(b$opt) # povprečje generiranih optimumov

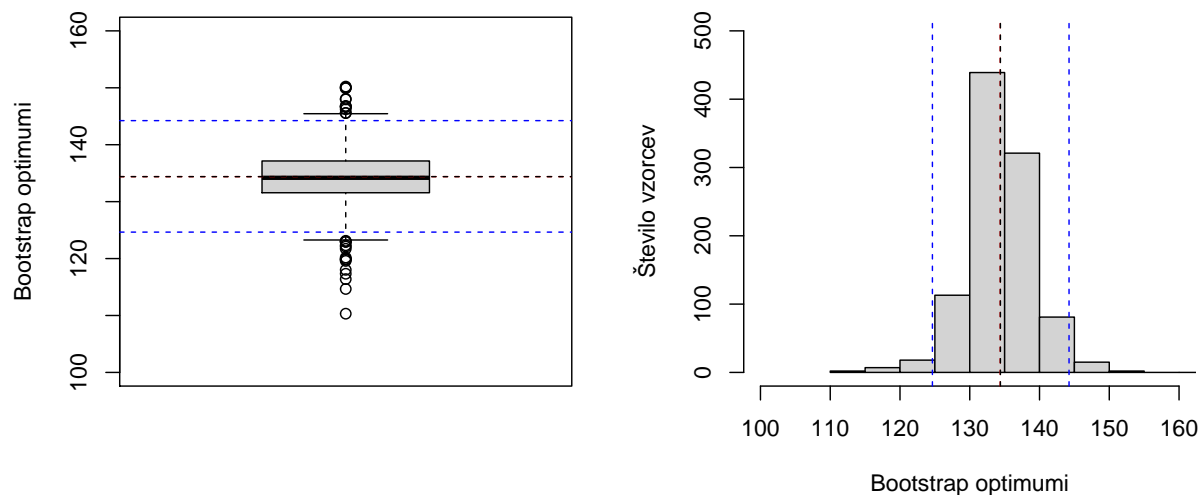
[1] 134.3672

> bias <- opt-mean(b$opt) # pristranskost bootstrap ocene
> (optIZ<-quantile(b$opt,c(0.025,0.975)) ) # centilni interval zaupanja za optimum

      2.5%      97.5%
124.6479 144.2324
```

Optimalna gostota setve je 134.4, pripadajoči 95 % centilni bootstrap IZ je (124.6, 144.2).

Grafični prikaz dobljene porazdelitve 1000 bootstrap optimumov je na Sliki 5.



Slika 5: Porazdelitev 1000 bootstrap optimumov za gostoto setve, okvir z ročaji (levo) in histogram (desno), črna črtkana črta predstavlja izračunani optimum na podlagi osnovnega modela, rdeča črtkana črta pa povprečje bootstrap optimumov, modre črtkane črte predstavljajo meje 95 % centilnega bootstrap intervala zaupanja za optimum

Za samovzorčenje primerov za `model.kvad` lahko uporabimo funkcijo `Boot` iz paketa `car`. Ta funkcija z osnovnimi argumenti naredi samovzorčenje primerov za neparametrični bootstrap za modele vrste `lm`, `glm` in `nls`.


```
> library(car)
> betahat.boot<-Boot(model.kvad, R=1000, f = coef, method = "case")
> summary(betahat.boot)
```

Number of bootstrap replications R = 1000

	original	bootBias	bootSE	bootMed
(Intercept)	0.24919269	-5.3285e-04	2.6288e-01	0.23967551
gostsetve	0.10358831	2.0743e-04	7.1856e-03	0.10387175
I(gostsetve^2)	-0.00038539	-2.2306e-06	3.9317e-05	-0.00038593

```
> head(betahat.boot$t)
```

	(Intercept)	gostsetve	I(gostsetve^2)
[1,]	0.316396257	0.10271458	-0.0003859367
[2,]	0.166486446	0.10269180	-0.0003725510
[3,]	0.270339981	0.10914402	-0.0004248780
[4,]	0.141740301	0.11224355	-0.0004278453
[5,]	0.785034593	0.08827226	-0.0003005972
[6,]	-0.002462006	0.10938614	-0.0004078204

```
> betahat.boot$t$opt<--0.5*betahat.boot$t[,2]/betahat.boot$t[,3]
> bootIZ<-quantile(betahat.boot$t$opt,c(0.025,0.975))
> round(cbind(mean(betahat.boot$t$opt),t(bootIZ)), 2)
```

	2.5%	97.5%
[1,]	134.38	145.02

Ali pa kar direkten izračun optimumov:

```
> # še druga možnost uporabe funkcije Boot za bootstrap optimume:
> set.seed(3435)
> betahat.boot.1<-Boot(object=model.kvad, R=1000, labels="optimum",
+                       f = function(object) -0.5*coef(object)[2]/coef(object)[3])
> summary(betahat.boot.1)
```

	R	original	bootBias	bootSE	bootMed
optimum 1000		134.39	0.25634	5.2841	134.52

```
> opt.bootstrap<-as.numeric(summary(betahat.boot.1)[2]+summary(betahat.boot.1)[3])
> round(opt.bootstrap,2)
```

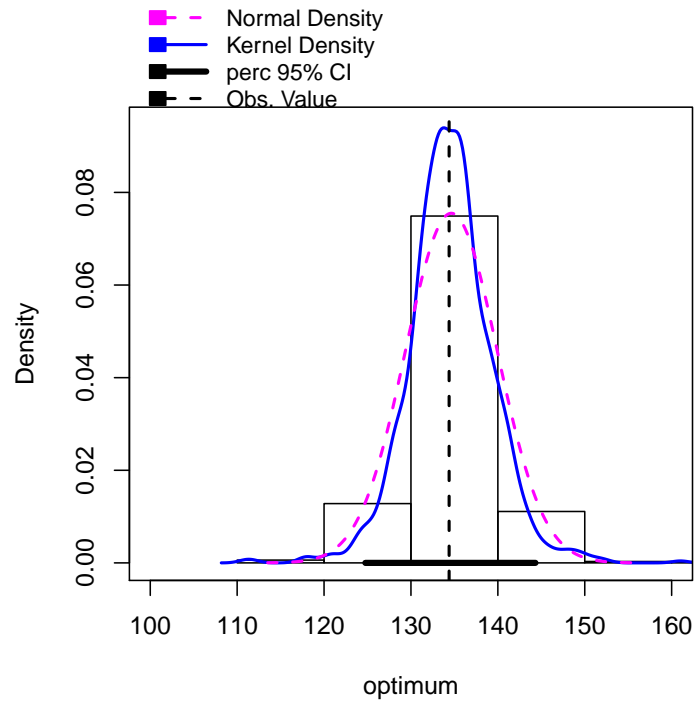
```
[1] 134.65
```

```
> confint(betahat.boot.1, type="perc")
```

Bootstrap percent confidence intervals

```
      2.5 %   97.5 %  
optimum 124.759 144.3246
```

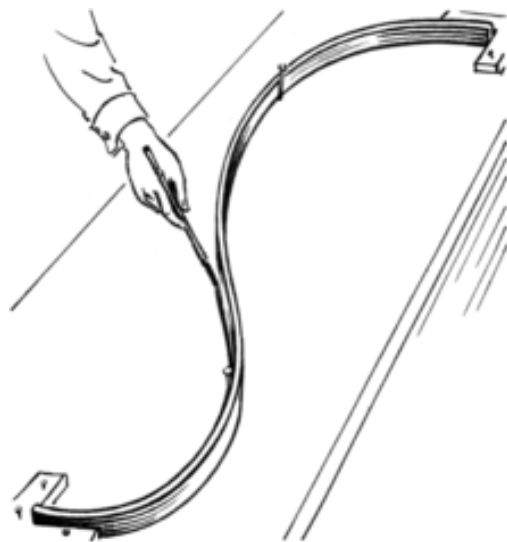
```
> hist(betahat.boot.1, ci="perc",, xlim=c(100, 160))
```



Slika 6: Histogram za bootstrap optimume gostote setve na podlagi model.kvad

1.2 Regresija zlepkov

Zlepek (*spline*) v angleškem jeziku predstavlja dolg in tenek upogljiv kos lesa ali kovine, ki so ga načrtovalci/konstruktorji uporabljali za risanje krivulj skozi vnaprej določene točke (Slika 7). Zlepke v regresijski analizi uporabljamo za opis nelinearnega odnosa med odzivno in izbrano napovedno spremenljivko.



Slika 7: Zlepek z dvema vozliščema oziroma s tremi odseki (Vir: Wikipedia)

Pogosto se zgodi, da nelinearnosti ne moremo opisati s polinomsko regresijo dovolj nizke stopnje. V takem primeru lahko vrednosti napovedne spremenljivke razdelimo na odseke in na posameznem odseku uporabimo polinomsko regresijo nižje stopnje. Takemu načinu modeliranja pravimo **polinomska regresija po odsekih** (*piecewise polynomials*). Na primer, če vrednosti spremenljivke x razdelimo na dva odseka: $x < c$ in $x \geq c$ in izberemo polinom tretje stopnje, v modelu polinomske regresije na odsekih ocenjujemo osem parametrov modela:

$$y_i = \begin{cases} \beta_{01} + \beta_{11}x_i + \beta_{21}x_i^2 + \beta_{31}x_i^3 + \varepsilon_i, & x_i < c, \\ \beta_{02} + \beta_{12}x_i + \beta_{22}x_i^2 + \beta_{32}x_i^3 + \varepsilon_i, & x_i \geq c. \end{cases} \quad (3)$$

Z modelom (3) podatkom prilagodimo dve polinomske funkciji, eno za podatke z $x_i < c$ in drugo za podatke z $x_i \geq c$. Vrednost spremenljivke $x = c$, kjer se vrednosti parametrov polinoma spremenijo, se imenuje **vozlišče** (*knot*). Več vozlišč omogoča bolj kompleksno nelinearno odvisnost. Če postavimo K vozlišč znotraj intervala vrednosti spremenljivke x , prilagodimo $K + 1$ polinomov izbrane stopnje p . V vozliščih je potrebno definirati, kako naj se polinoma stikata. Najbolj uporabno je, da se polinoma stikata zvezno in gladko, v takem primeru govorimo o **regresiji zlepkov**.

Zlepek je funkcija, ki opiše krivuljo na izbranih odsekih napovedne spremenljivke x . Odseki so določeni z vozlišči. Na posameznem odseku odvisnost odzivne spremenljivke od x opiše

polinom p -te stopnje. V vozliščih se vrednosti sosednjih polinomov gladko stikajo, kar pomeni, da pri ocenjevanju parametrov polinomov postavimo še dodatne pogoje: v vozlišču morata imeti sosednja polinoma stopnje p isto vrednost in iste vrednosti odvodov reda od $1, \dots, p-1$. Ti dodatni pogoji zmanjšajo število parametrov, ki jih moramo oceniti v modelu.

Število vozlišč K izberemo vnaprej, prav tako njihove vrednosti; izbira je odvisna od števila podatkov, kompleksnosti nelinearnosti na posameznih odsekih in od predhodnega poznavanja procesa, ki ga modeliramo.

1.2.1 Bazne funkcije

Za razumevanje regresije zlepkov najprej definirajmo t. i. **bazne funkcije**. Polinomska regresija predstavlja poseben primer pristopa regresijskega modeliranja z baznimi funkcijami. Ideja baznih funkcij je v tem, da napovedno spremenljivko x v model vključimo v obliki različnih transformacij oziroma v obliki k baznih funkcij: $b_1(x), b_2(x), \dots, b_k(x)$. V modelu **bazne funkcije predstavljajo regresorje**:

$$y_i = \beta_0 + \beta_1 b_1(x_i) + \beta_2 b_2(x_i) + \beta_3 b_3(x_i) + \dots + \beta_k b_k(x_i) + \varepsilon_i. \quad (4)$$

Bazne funkcije so vedno izbrane vnaprej. V kontekstu polinomske regresije stopnje p imamo $k = p$ baznih funkcij $b_j(x) = x^j$, $j = 1, \dots, p$. Če gre za ortogonalne polinome, potem bazne funkcije predstavljajo linearno kombinacijo regresorjev x^j , $j = 1, \dots, p$, ki ustreza pogoju, da so bazne funkcije med seboj neodvisne.

V primeru modeliranja z baznimi funkcijami parametre modela (4) ocenimo po metodi najmanjših kvadratov. Če so splošne predpostavke linearnega modela izpolnjene, je inferenca na ocenah parametrov enaka kot v primeru navadnih regresorjev.

Bazne funkcije lahko predstavljajo zelo različne funkcije napovedne spremenljivke, zelo pogosto so določene kot kombinacija polinomov nižjih stopenj (največ tretje).

1.2.2 Linearni zlepki

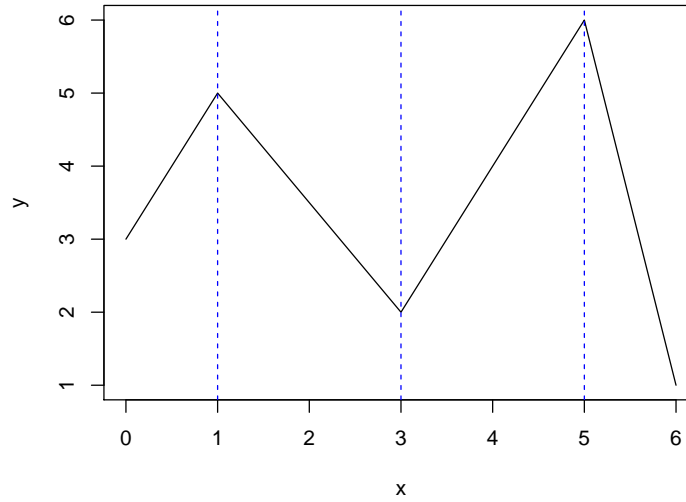
Linearni zlepki predstavljajo linearno funkcijo, ki se lomi v K vozliščih. Za primer pogledjmo linearne zlepke s tremi vozlišči $K = 3$. Vrednosti napovedne spremenljivke x so razdeljene na štiri odseke pri vozliščih $x = a_1$, $x = a_2$ in $x = a_3$. Model linearnega zlepkov predstavlja funkcija:

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 (x_i - a_1)_+ + \beta_3 (x_i - a_2)_+ + \beta_4 (x_i - a_3)_+ + \varepsilon_i, \quad (5)$$

kjer velja $(u)_+ = u$, $u > 0$ in $(u)_+ = 0$, $u \leq 0$. V (5) je x_i bazna funkcija, $(x_i - a_1)_+$, $(x_i - a_2)_+$ in $(x_i - a_3)_+$ pa so t. i. **odrezane bazne funkcije** (*truncated basis*).

Enačbo (5) lahko zapišemo po odsekih napovedne spremenljivke x :

$$\begin{aligned}
 y_i &= \beta_0 + \beta_1 x_i + \varepsilon_i, & x_i &\leq a, \\
 \beta_0 + \beta_1 x_i + \beta_2(x_i - a_1) + \varepsilon_i, & & a_1 < x_i &\leq a_2, \\
 \beta_0 + \beta_1 x_i + \beta_2(x_i - a_1) + \beta_3(x_i - a_2) + \varepsilon_i, & & a_2 < x_i &\leq a_3, \\
 \beta_0 + \beta_1 x_i + \beta_2(x_i - a_1) + \beta_3(x_i - a_2) + \beta_4(x_i - a_3) + \varepsilon_i, & & a_3 < x_i. & \quad (6)
 \end{aligned}$$



Slika 8: Linearni zlepek z vozlišči pri $a_1 = 1$, $a_2 = 3$ in $a_3 = 5$

Linearne zlepke v regresijski model s K vozlišči pri vrednostih (a_1, a_2, \dots, a_K) vključimo z baznimi funkcijami $b_1(x) = x$, $b_2(x) = (x - a_1)_+$ do $b_{K+1}(x) = (x - a_K)_+$, njihovo število je določeno s številom vozišč $K + 1$:

$$y_i = \beta_0 + \beta_1 b_1(x_i) + \beta_2 b_2(x_i) + \dots + \beta_{K+1} b_{K+1}(x_i) + \varepsilon_i. \quad (7)$$

Ocene $K + 1$ parametrov modela (7) izračunamo po metodi najmanjših kvadratov ob dodatnem pogoju, da se vrednosti \hat{y} stikajo v vozliščih. Linearnost odvisnosti y od x testiramo z ničelno domnevo $H_0 : \beta_2 = \beta_3 = \dots = \beta_{K+1} = 0$. Uporabimo F -test za dva gnezdena modela.

Linearni zleпки so preprosti in z njimi lahko opišemo veliko odnosov, njihova slabost pa je, da se funkcija v vozliščih prelomi. Če želimo modelirati gladke krivulje, moramo za opis nelinearnosti na posameznih odsekih uporabiti polinome višjih stopenj.

1.2.3 Kubični zlepki

Praksa je pokazala, da imajo polinomi tretje stopnje (kubični polinomi) lepe lastnosti in sposobnost, da ob primerni izbiri števila vozlišč opišejo tudi zelo kompleksne nelinearne odvisnosti. Dva kubična polinoma se gladko stikata v vozlišču, če v vozlišču poleg vrednosti izenačimo tudi njun prvi in drugi odvod. Model **kubičnega zleпка** za tri vozlišča ($K = 3$) je:

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \beta_3 x_i^3 + \beta_4 (x_i - a_1)_+^3 + \beta_5 (x_i - a_2)_+^3 + \beta_6 (x_i - a_3)_+^3 + \varepsilon_i, \quad (8)$$

kjer so bazne funkcije in odrezane potenčne bazne funkcije (*truncated power basis*)

$$\begin{aligned} b_1(x) &= x, & b_2(x) &= x^2, & b_3(x) &= x^3, \\ b_4(x) &= (x - a_1)_+^3, & b_5(x) &= (x - a_2)_+^3, & b_6(x) &= (x - a_3)_+^3. \end{aligned} \quad (9)$$

Model kubičnega zleпка za K -vozlišč je:

$$y_i = \beta_0 + \beta_1 b_1(x_i) + \beta_2 b_2(x_i) + \dots + \beta_{K+3} b_{K+3}(x_i) + \varepsilon_i, \quad (10)$$

kjer so prve tri bazne funkcije določene enako kot v (9), vse naslednje pa:

$$b(x, a_j) = (x - a_j)_+^3 = \begin{cases} (x - a_j)^3, & x > a_j \\ 0, & \text{drugače,} \end{cases} \quad (11)$$

$a_j, j = 1, \dots, K$ so vozlišča.

Parametri modela regresijskih zlepkov ($\beta_0, \dots, \beta_{K+3}$) so izračunani po metodi najmanjših kvadratov z upoštevanimi dodatnimi pogoji, ki zagotavljajo, da so njihovi stiki v vozliščih gladki. Če ima kubični zlepek K vozlišč, moramo v regresijskem modelu poleg presečišča oceniti $K + 3$ parametrov.

1.2.4 Naravni zlepki

Praksa je pokazala, da se pri kubičnih regresijskih zlepkih pogosto zgodi, da se slabo obnesejo na prvem in zadnjem odseku (pred prvim in za zadnjim vozliščem). To težavo rešimo z uporabo t. i. **naravnih zlepkov** (*natural splines* ali *restricted cubic splines*). V tem primeru predpostavimo linearni odnos med y in x na prvem in zadnjem odseku. Posledično v modelu ocenjujemo poleg presečišča samo $K - 1$ parametrov: odpadeta parametra pri x^2 in x^3 , zadnja dva parametra β_K in β_{K+1} se zapišeta kot linearna kombinacija predhodnih parametrov $\beta_2, \dots, \beta_{K-1}$ (16). Regresijski model z naravnimi zlepkami zapišemo:

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 (x_i - a_1)_+^3 + \beta_3 (x_i - a_2)_+^3 + \dots + \beta_{K+1} (x_i - a_K)_+^3 + \varepsilon, \quad (12)$$

a_1, \dots, a_k so vozlišča. V procesu ocenjevanja parametrov modela naravnih zlepkov najprej ocenjujemo K parametrov na podlagi funkcije:

$$y_i = \beta_0 + \beta_1 b_1(x_i) + \beta_2 b_2(x_i) + \dots + \beta_{K-1} b_{K-1}(x_i), \quad (13)$$

kjer je $b_1(x_i) = x_i$, ostale bazne funkcije so za $j = 2, \dots, K-1$ izražene takole:

$$b_{j+1}(x_i) = (x - a_j)_+^3 - (x - a_{K-1})_+^3 (a_K - t_j) / (t_K - a_{K-1}) + (x - a_K)_+^3 (a_{K-1} - a_j) / (a_K - a_{K-1}). \quad (14)$$

Pokažemo lahko, da je bazna funkcija $b_{K+1}(x)$ na zadnjem odseku ($x \geq a_K$) linearna. Na podlagi ocen $\hat{\beta}_0, \dots, \hat{\beta}_{K-1}$ se izračuna še oceni $\hat{\beta}_K$ in $\hat{\beta}_{K+1}$ iz (12):

$$\hat{\beta}_K = [\hat{\beta}_2(a_1 - a_K) + \hat{\beta}_3(a_2 - a_K) + \dots + \hat{\beta}_{K-1}(a_{K-2} - a_K)] / (a_K - a_{K-1}), \quad (15)$$

$$\hat{\beta}_{K+1} = [\hat{\beta}_2(a_1 - a_{K-1}) + \hat{\beta}_3(a_2 - a_{K-1}) + \dots + \hat{\beta}_{K-1}(a_{K-2} - a_{K-1})] / (a_{K-1} - a_K). \quad (16)$$

Pri modeliranju regresijskih zlepkov je število in položaj vozlišč določeno vnaprej. Položaj vozlišč lahko določimo na podlagi predhodnega poznavanja procesa, ki ga modeliramo; na primer, če vemo, da se naklon spremeni pri vrednosti $x = a$, vrednost a vnaprej izberemo za vozlišče. V splošnem smiselni vrednosti za vozlišča ne poznamo. Analize so pokazale, da samo položaj vozlišč ni tako pomemben, bolj pomembno je število vozlišč. Položaj vozlišč po navadi določimo z vrednostmi enako razmaknjenih kvantilov napovedne spremenljivke x . S tem zagotovimo, da je število podatkov na vseh odsekih uravnoteženo. Pogosto sta prvi in zadnji odsek ob uporabi naravnih zlepkov manjša, priporočene kvantile kaže Tabela 1.

Tabela 1: Priporočeni kvantilni rangi za vozlišča naravnih zlepkov (Harrell F. E., 2015)

Število vozlišč k	Kvantilni rangi
3	.10, .5, .90
4	.05, .35, .65, .95
5	.05, .275, .5, .725, .95
6	.05, .23, .41, .59, .77, .95
7	.025, .1833, .3417, .5, .6583, .8167, .975

Če imamo malo podatkov ($n \leq 30$), običajno izberemo $K = 3$, sicer je najpogostejša primerna izbira $K = 4$ ali $K = 5$. Za velik n ($n \geq 100$), je običajno primerno število vozlišč $K = 7$, večje vrednosti za K so zelo redko potrebne.

Število potrebnih vozlišč v praksi pogosto določimo na podlagi navzkrižnega preverjanja modela (*cross-validation*). Ta postopek bomo spoznali v enem izmed poglavij, ki sledijo.

1.2.5 Primer: KORUZA (nadaljevanje)

Regresijo naravnih zlepkov bomo najprej uporabili na primeru napovedi pridelka koruze v odvisnosti od gostote setve (datoteka KORUZA.txt). Za regresijo zlepkov potrebujemo paket `splines`. Za modeliranje zlepkov stopnje p s K vozlišči uporabimo funkcijo `bs`, za modeliranje naravnih zlepkov pa funkcijo `ns`.

Najprej pogledjmo argumente funkcije `bs(x, df = NULL, knots = NULL, degree = 3, intercept = FALSE, Boundary.knots = range(x))`:

- prvi argument funkcije `bs` je vektor napovedne spremenljivke x , v našem primeru bo to gostota setve (`gostsetve`). Če je to edini argument, funkcija vrne tri bazne funkcije za polinom tretje stopnje, ker ima po prednastavitvi argument `degree` vrednost 3 ($p = 3$), argument `knots` pa `NULL` ($K = 0$);

```
> library(splines)
> # df=NULL, knots=NULL
> bs.0<-bs(koruza$gostsetve)
> str(bs.0)

'bs' num [1:45, 1:3] 0.442 0.264 0.143 0.433 0.425 ...
- attr(*, "dimnames")=List of 2
..$ : NULL
..$ : chr [1:3] "1" "2" "3"
- attr(*, "degree")= int 3
- attr(*, "knots")= num(0)
- attr(*, "Boundary.knots")= num [1:2] 4.97 172.1
- attr(*, "intercept")= logi FALSE

> head(bs.0)

           1           2           3
[1,] 0.4422797 0.24939741 0.046877627
[2,] 0.2641465 0.03316373 0.001387908
[3,] 0.1429673 0.42325434 0.417681154
[4,] 0.4325936 0.16188647 0.020193880
[5,] 0.4247011 0.14552358 0.016621190
[6,] 0.0000000 0.00000000 0.000000000
```

- argument `df` predstavlja število stopinj prostosti regresijskega modela `degree+K`; s tem argumentom lahko posredno nastavimo število vozlišč $K = df - degree$; po prednastavitvi ima vrednost `NULL`, kar pomeni, da je $K = 0$;

```
> # df=4, knots=NULL
> bs.1<-bs(koruza$gostsetve, df=4)
> str(bs.1)
```



```
'bs' num [1:45, 1:4] 0.4089 0.5816 0.0252 0.5797 0.6071 ...
- attr(*, "dimnames")=List of 2
..$ : NULL
..$ : chr [1:4] "1" "2" "3" "4"
- attr(*, "degree")= int 3
- attr(*, "knots")= Named num 65.2
..- attr(*, "names")= chr "50%"
- attr(*, "Boundary.knots")= num [1:2] 4.97 172.1
- attr(*, "intercept")= logi FALSE
```

```
> head(bs.1)
```

```
      1      2      3      4
[1,] 0.40887184 0.46111806 0.130010096 0.00000000
[2,] 0.58157783 0.08514991 0.003849216 0.00000000
[3,] 0.02517428 0.20938968 0.543850631 0.2215854
[4,] 0.57967645 0.34965495 0.056005570 0.00000000
[5,] 0.60710389 0.32184584 0.046097098 0.00000000
[6,] 0.00000000 0.00000000 0.000000000 0.00000000
```

- z argumentom **knots** nastavimo položaje vozlišč, vrednosti vozlišč zapišemo v vektor. Če ima vrednost NULL in je argument **df** različen od NULL, se položaji vozlišč določijo na podlagi kvantilnih rangov, ki vrednosti za x razdelijo na enake dele glede na število vozlišč. Na primer če je **df**=5 in **degree**=3, vozlišča predstavljata 33.3 centil in 66.7. centil vrednosti x ;
- argument **intercept** ima po prednastavitvi vrednost **FALSE**, kar pomeni, da presečišče ni vključeno pri računanju baznih funkcij zleпка, to je priročno za uporabo funkcije **bs** v formuli modela **lm**;
- argument **Boundary.knots** ima po prednastavitvi vrednosti **min(x)** in **max(x)** ter določa razpon vrednosti spremenljivke x , na katerem se računajo bazne funkcije zlepkov.

Ilustracija baznih funkcij kubičnih zlepkov s tremi vozlišči določenimi s kvartili **gostsetve**:

```
> # določimo vrednosti gostsetve za vozlišča
> voz1<-quantile(koruza$gostsetve, probs = c(0.25, 0.5, 0.75), na.rm=T)
> bs.3<-bs(koruza$gostsetve, knots=voz1, degree=3)
> # enakovreden ukaz je
> # bs.3<-bs(koruza$gostsetve,df=6)
> str(bs.3)

'bs' num [1:45, 1:6] 0 0.519 0 0.0487 0.0829 ...
- attr(*, "dimnames")=List of 2
..$ : NULL
..$ : chr [1:6] "1" "2" "3" "4" ...
```

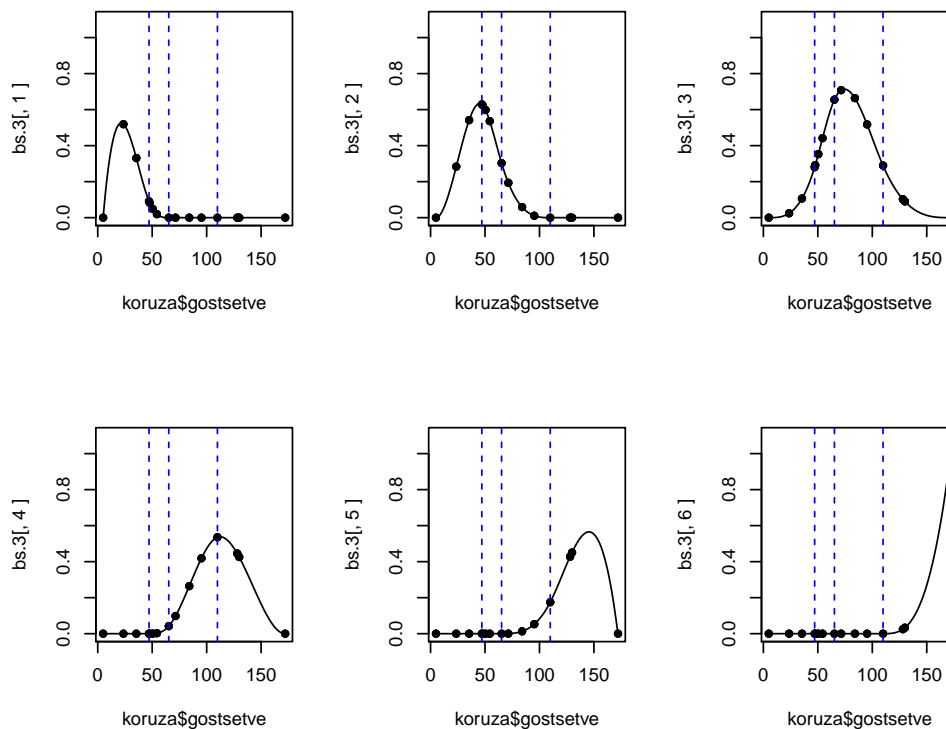
```
- attr(*, "degree")= int 3
- attr(*, "knots")= Named num [1:3] 47.1 65.2 109.9
  ..- attr(*, "names")= chr [1:3] "25%" "50%" "75%"
- attr(*, "Boundary.knots")= num [1:2] 4.97 172.1
- attr(*, "intercept")= logi FALSE

> # vrednosti baznih funkcij zlepkov polinomov tretje stopnje s tremi vozlišči
> # za dane vrednosti gostsetve
> head(bs.3)
```

	1	2	3	4	5	6
[1,]	0.00000000	0.3027066	0.65534883	4.194458e-02	0.0000000	0.00000000
[2,]	0.51904398	0.2835271	0.02433147	0.000000e+00	0.0000000	0.00000000
[3,]	0.00000000	0.0000000	0.09042982	4.257325e-01	0.4505933	0.03324443
[4,]	0.04868531	0.5981136	0.35292537	2.757292e-04	0.0000000	0.00000000
[5,]	0.08285142	0.6257652	0.29138232	1.104656e-06	0.0000000	0.00000000
[6,]	0.00000000	0.0000000	0.00000000	0.000000e+00	0.0000000	0.00000000

Za ilustracijo Slika 9 prikazuje vrednosti baznih funkcij za regresijo kubičnih zlepkov s tremi vozlišči.

```
> par(mfrow=c(2,dim(bs.3)[2]/2))
> x<-seq(min(koruza$gostsetve), max(koruza$gostsetve),1 )
> for (i in 1:dim(bs.3)[2])
+ {plot(koruza$gostsetve,bs.3[,i], pch=16, ylim=c(0,1.1),
+       ylab=paste("bs.3[,", i,""])))
+   lines(x, bs(x, knots=voz1)[,i])
+   abline(v=voz1, col="blue", lty=2)}
```



Slika 9: Grafična predstavitev baznih funkcij kubičnega zleпка s tremi vozlišči, $K = 3$, ki predstavljajo kvartile gostsetve

Modelirajmo odvisnost pridelka koruze od gostote setve z uporabo zlepkov. Najprej bomo naredili model, ki je enakovreden modelu polinomske regresije reda 2. Funkcija `bs` ima argument `degree = 2`, število vozlišč je 0 (`knots=NULL`, prednastavitev):

```
> # model regresije kvadratnih zlepkov brez vozlišč
> model.bs2.0<-lm(prid1.ha ~ bs(gostsetve, degree=2), data=koruza)
> coef(summary(model.bs2.0))
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	0.7542077	0.2858920	2.638086	1.164340e-02

```
bs(gostsetve, degree = 2)1 8.3365769 0.6603105 12.625238 6.903840e-16
bs(gostsetve, degree = 2)2 5.9077512 0.3849820 15.345524 7.976758e-19
```

```
> # za primerjavo izpišimo ocene parametrov polinomske regresije
> coef(summary(model.kvad))
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	0.2491926865	3.218808e-01	0.774177	4.431625e-01
gostsetve	0.1035883074	8.341093e-03	12.419033	1.195171e-15
I(gostsetve^2)	-0.0003853948	4.612724e-05	-8.355037	1.797520e-10

```
> # še polinomska regresija z baznimi funkcijami,
> # ki predstavljajo ortogonalne kvadratne polinome regresorjev, degree=2, raw=FALSE
> model.kvad.1 <-lm(prid1.ha ~ poly(gostsetve, 2), data=koruza)
> coef(summary(model.kvad.1))
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	5.095022	0.1006226	50.634946	2.781172e-39
poly(gostsetve, 2)1	10.746839	0.6749972	15.921309	2.119810e-19
poly(gostsetve, 2)2	-5.639627	0.6749972	-8.355037	1.797520e-10

```
> # koeficienti determinacije za vse tri modele
> summary(model.kvad)$r.squared
```

```
[1] 0.8850243
```

```
> summary(model.kvad.1)$r.squared
```

```
[1] 0.8850243
```

```
> summary(model.bs2.0)$r.squared
```

```
[1] 0.8850243
```

Ocene parametrov za `model.kvad`, `model.kvad.1` in `model.bs2.0` so različne, ker so model-ske matrike različne, koeficienti determinacije pa so isti in skoraj identične so tudi napovedi modelov (Slika 10).

```
> head(model.matrix(model.kvad))
```

	(Intercept)	gostsetve	I(gostsetve^2)
1	1	65.230	4254.95290
2	1	23.610	557.43210
3	1	129.900	16874.01000
4	1	50.480	2548.23040
5	1	47.620	2267.66440
6	1	4.967	24.67109

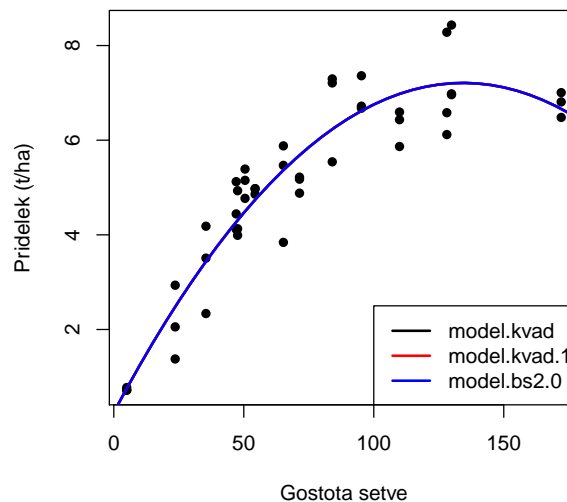
```
> head(model.matrix(model.kvad.1))
```

	(Intercept)	poly(gostsetve, 2)1	poly(gostsetve, 2)2
1	1	-0.03201627	-0.10913359
2	1	-0.17374630	0.13258280
3	1	0.18820671	-0.01498477
4	1	-0.08224496	-0.05055419
5	1	-0.09198421	-0.03575398
6	1	-0.23723195	0.31763110

```
> head(model.matrix(model.bs2.0))
```

	(Intercept)	bs(gostsetve, degree = 2)1	bs(gostsetve, degree = 2)2
1	1	0.4611181	0.13001010
2	1	0.1982068	0.01244249
3	1	0.3774811	0.55876593
4	1	0.3963200	0.07415604
5	1	0.3801498	0.06512905
6	1	0.0000000	0.00000000

```
> novi.x<-data.frame(gostsetve=seq(0, 180, 5))
> plot(koruza$gostsetve,koruza$prid1.ha, pch=16,
+      ylab="Pridelek (t/ha)", xlab="Gostota setve",)
> lines(novi.x$gostsetve, predict(model.kvad, novi.x), lwd=2, col="black")
> lines(novi.x$gostsetve, predict(model.kvad.1, novi.x),lwd=2, col="red")
> lines(novi.x$gostsetve, predict(model.bs2.0, novi.x),lwd=2, col="blue")
> legend(100, 2.5, legend=c("model.kvad","model.kvad.1","model.bs2.0"),
+       col=c("black","red","blue"), lwd=2, lty=1)
```



Slika 10: Odvisnost pridelka (t/ha) od gostote setve: parabola izračunana z `model.kvad` in napovedi za `model.kvad.1` in `model.bs2.0`

Modelov `model.bs2.0` in `model.kvad` se ne da primerjati z F -testom, ker modela nista gnezdena, imata enako število parametrov. Vsoti kvadriranih ostankov sta enaki:

```
> anova(model.kvad, model.bs2.0)
```

Analysis of Variance Table

Model 1: prid1.ha ~ gostsetve + I(gostsetve^2)

Model 2: prid1.ha ~ bs(gostsetve, degree = 2)

	Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
1	42	19.136				
2	42	19.136	0	-7.1054e-15		

V drugem primeru bomo pridelek koruze modelirali s kubičnimi in z naravnimi zlepkami s tremi vozlišči določenimi s kvantilnimi rangi 0.25, 0.5 in 0.75.

```
> model.bs.3<-lm(prid1.ha ~ bs(gostsetve, knots=voz1), data=koruza)
> model.ns.3<-lm(prid1.ha ~ ns(gostsetve, knots=voz1), data=koruza)
> coef(summary(model.bs.3))
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	0.7384488	0.3964574	1.8626184	7.026053e-02
bs(gostsetve, knots = voz1)1	0.3605557	0.9693291	0.3719642	7.119858e-01
bs(gostsetve, knots = voz1)2	3.7417868	0.6514026	5.7442000	1.277304e-06
bs(gostsetve, knots = voz1)3	4.8728190	0.6888817	7.0735212	1.956330e-08
bs(gostsetve, knots = voz1)4	6.8665761	1.1516307	5.9624808	6.409079e-07
bs(gostsetve, knots = voz1)5	6.3003290	1.3579431	4.6396119	4.064820e-05
bs(gostsetve, knots = voz1)6	6.0334575	0.5611930	10.7511268	4.391166e-13

```
> coef(summary(model.ns.3))
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	0.5674518	0.3660588	1.550166	1.289785e-01
ns(gostsetve, knots = voz1)1	5.2579716	0.4883474	10.766867	2.197480e-13
ns(gostsetve, knots = voz1)2	5.7707941	0.4793138	12.039701	7.050328e-15
ns(gostsetve, knots = voz1)3	9.5461236	0.8772115	10.882351	1.594601e-13
ns(gostsetve, knots = voz1)4	4.2841747	0.4499559	9.521322	7.793808e-12

Primerjava modela polinomske regresije in regresijskih zlepkov:

```
> anova(model.kvad, model.bs.3)
```

Analysis of Variance Table

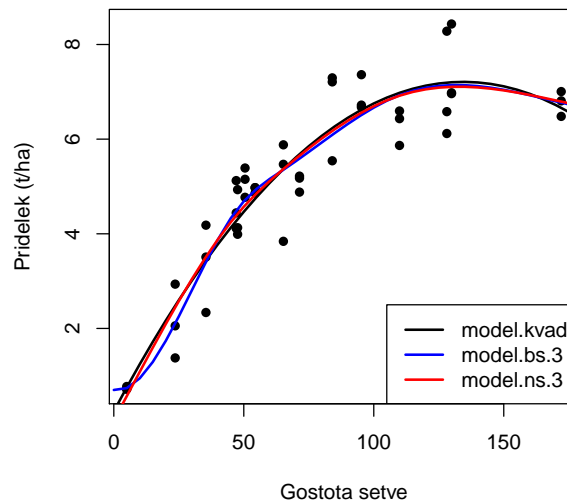
```
Model 1: prid1.ha ~ gostsetve + I(gostsetve^2)
Model 2: prid1.ha ~ bs(gostsetve, knots = voz1)
  Res.Df    RSS Df Sum of Sq    F Pr(>F)
1      42 19.136
2      38 17.988  4    1.1478 0.6062 0.6606
```

```
> anova(model.kvad, model.ns.3)
```

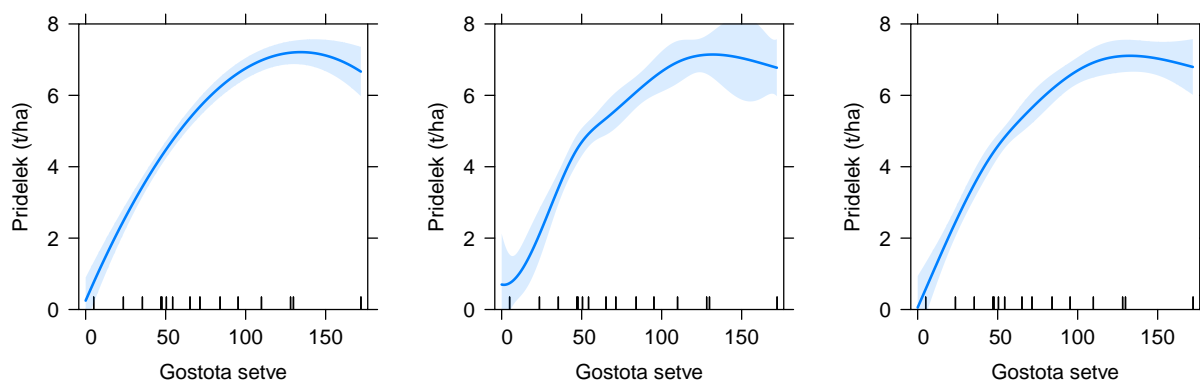
Analysis of Variance Table

```
Model 1: prid1.ha ~ gostsetve + I(gostsetve^2)
Model 2: prid1.ha ~ ns(gostsetve, knots = voz1)
  Res.Df    RSS Df Sum of Sq    F Pr(>F)
1      42 19.136
2      40 18.645  2    0.49135 0.5271 0.5944
```

Za `model.kvad` velja, da se pridelek od optimalne gostote naprej enako hitro zmanjšuje, kot se je povečeval pred dosegom optimuma. Pri modelih z regresijskimi zlepkmi `model.bs.3` in `model.ns.3` pa je padec pridelka po optimumu počasnejši (Sliki 11, 12). Med modeli ni statistično značilnih razlik v pojasnjeni variabilnosti odzivne spremenljivke.



Slika 11: Odvisnost pridelka (t/ha) od gostote setve: parabola izračunana z `model.kvad` in napovedi za `model.bs.3` in `model.ns.3`



Slika 12: Odvisnost pridelka (t/ha) od gostote setve, napovedi za `model.kvad` (levo), `model.bs.3` (sredina) in za `model.ns.3` (desno) ter 95 % intervali zaupanja za povprečne napovedi

2 VAJE

2.1 Telesna masa in višina žensk

Analizirajte odvisnost telesne mase od telesne višine za ženske stare med 30 in 39 let. Podatki so v podatkovnem okviru `women` v paketu `stats`. Podatke grafično prikažite, naredite ustrezen model, preverite predpostavke izbranega modela, napovedi grafično prikažite in napišite obrazložitev statistične analize.

2.2 Plača

V podatkovnem okviru `Wage` v paketu `ISLR` so podatki o plačah 3000 moških delavcev v srednje atlantski regiji. Analizirajmo odvisnost plače (`wage`) od starosti (`age`), leta pridobitve podatkov (`year`) in o izobrazbi delavcev (`education`).