

Datorlaboration 2

Josef Wilzén

27 augusti 2025

732G57 Maskininlärning för statistiker HT2025

Allmänt

Datorlaborationerna kräver att ni har R och Rstudio installerat.

- Kodmanual: [länk](#)
- **ISL**: An introduction to Statistical Learning,
 - Boken: [länk](#)
 - R-kod till labbar: [länk](#)
 - Dataset: [länko](#)ch [länk](#)
- **IDM**: Introduction to Data Mining,
 - Kod till boken: [länk](#)
 - [Sample chapters](#)
- Dataset till vissa uppgifter finns [här](#)
- Funktionerna `lm_diagnostics()` och `model_diagnostics()` ger utvärderingsplottar för residuallerna. Dessa finns här: [länk](#)

Notera att ni inte behöver göra alla delar på alla uppgifter. Det viktiga är att ni får en förståelse för de olika principerna och modellerna som avhandlats. Dessa uppgifter ska inte lämnas in, utan är till för er övning.

Datoruppdelning

För att motverka överanpassning bör ni dela upp data till träning-, validering-, (och testmängd). Detta kan göras med `createDataPartition()` från `caret`-paketet. Argument till den funktionen som är av vikt här är `p` som anger hur stor andel av observationerna som ska användas till träningsmängden. Ni kan också använda `subset()` för att göra detta, men det blir svårare att tydligt ange de observationer som ska tilldelas till valideringsmängden. Denna uppdelning måste var aslumpmässigt. Notera att om en testmängd ska skapas måste uppdelningen ske en gång till från valideringsmängden.

1 Polynomregression, stegfunktioner och basfunktioner

1.1 Introduktion

Det finns många sätt som vi kan utöka den vanliga linjära regressionsmodellen för att modellera icke-linjära samband mellan de förklarande variablerna och responsvariabeln. Notera att detta är ett stort fält med en rik litteratur, med många ofta avancerade regressionsmodeller. Vi ska här undersöka några enkla utökningar av linjär regression.

Vi kan beskriva regressionsmodeller som

$$y_i = f(x_i) + E \quad E \sim N(0, \sigma^2)$$

där $f(x)$ är en funktion som modellerar sambandet mellan x och y .

I fallet med den linjära regressionsmodellen:

$$f(x_i) = \beta_0 + \beta_1 x_i \Leftrightarrow y_i = \beta_0 + \beta_1 x_i + E$$

I detta fall så är $f(x)$ en *linjär* funktion. Vi kan såklart låta $f()$ bero av fler variabler:

$$f(x_{i,1}, x_{i,2}) = \beta_0 + \beta_1 x_{i,1} + \beta_2 x_{i,2}$$

eller mer generellt

$$f(\mathbf{x}_i) = \mathbf{x}_i^T \boldsymbol{\beta}$$

där $\mathbf{x}^T = (1, x_{i,1}, x_{i,2}, \dots, x_{i,p})$ och $\boldsymbol{\beta}^T = (\beta_0, \beta_1, \beta_2, \dots, \beta_p)$.

1. Vi utgår från att vi har *en* förklarande variabel här, x , men många av de saker som presenteras nedan går att relativt enkelt utökas till flera förklarande variabler.
2. För att undersöka vad för typ av samband som finns mellan x och y så är det oftast lättast plotta vår data i en scatter plot. Då kan se vad för typ av samband det finns mellan x och y . Ibland så kan vi göra lämpliga transformationer av y för att få sambanden mellan x och y mer linjärt, tex logaritmtransformation¹.
 - (a) Är de oberoende? Dvs värdena på y ändras inte när värden på x ändras.
 - (b) Linjärt beroende: y ökar eller minskar linjärt när värden på x ökar.
 - (c) Icke-linjärt beroende: Här kan vi ha många olika varianter på hur y ändras när x ändras. Exempel på icke-linjära samband är:

$$y = x^2 + E \quad y = \sin(x) + E \quad y = \tanh(x) + 0.1 \cdot x^3 + E$$

3. Om vi har skäl att tro att x och y har ett icke-linjärt samband, så kan vi skapa en eller flera nya variabler genom att använda icke-linjära transformationer på x . Ett vanligt exempel är en polynomtransformation²: x^2, x^3, \dots, x^h Sen använder vi dessa som vanliga kovariater i en linjär regressionsmodell. Vi får då att

$$f(x) = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \dots + \beta_h x^h \quad (1)$$

Vi kan spara dessa transformationer i nya variabler $z_1 = x$, $z_2 = x^2$, $z_3 = x^3$ osv, och sen använder vi dessa i vår designmatris $X = (1, z_1, z_2, \dots, z_h)$, och modellen blir då $y = X\boldsymbol{\beta} + E$.

4. Vi antar då att y beror linjärt på de nya variablerna z_1 till z_h , och vi kan skatta $\boldsymbol{\beta}$ med OLS. Så vårt modellantagande är att vi antar att y beror linjärt på de nya variablerna, men att y beror *icke-linjärt* på de ursprungliga variabeln x . Så denna idé att vi har variabler, som vi sen transformerar icke-linjärt och som vi sen använder som de nya variablerna i vår designmatris kan vi generalisera till andra transformationer än polynomtransformationer.
5. Notera: Polynom av av tillräckligt hög ordning kan anpassa kontinuerliga funktioner godtyckligt bra. Så om vi har tillräckligt högt h i funktionen i (1) kan vi i många fall få en bra anpassning av y . Dock så ökar risken för överanpassning ju högre h vi har. Det finns vissa problem med höga ordningens polynom (känsliga för outliers tex), ofta så väljer vi tredje eller fjärdegradspolynom som högst.
6. Vi låter $z = g(x)$ vara en funktion som transformerar x . Exempel:

- (a) $z = g(x) = \log(x)$
- (b) $z = g(x) = \exp(x)$
- (c) $z = g(x) = \sin(x)$
- (d) $z = g(x) = \tanh(x)$

¹Notera: $\log()$ syftar på logartimen med bas e i laborationerna

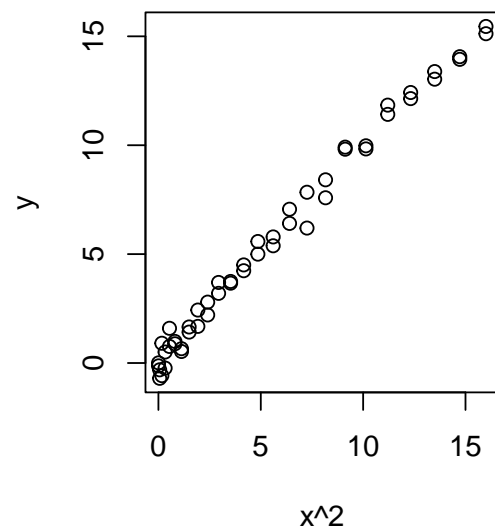
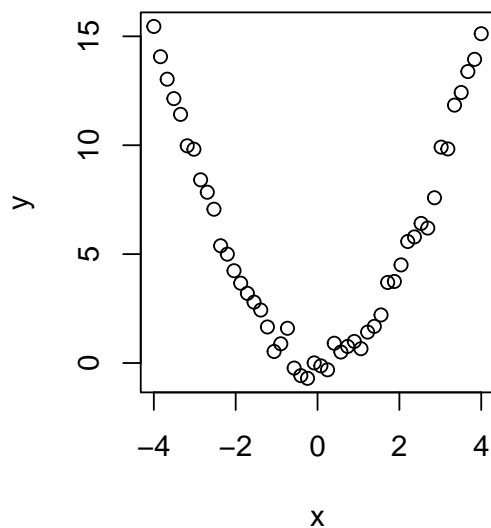
²Notera: skapar vi polynom av x så vill vi i regel centrera eller standardisera x först.

(e) $z = g(x) = \sqrt{x}$

(f) $z = g(x) = 0.5 \cdot \log(x) + 2 \cdot \sin(x)$

7. Se ett exempel nedan. Sambandet mellan x och y är kvadratisk, men mellan x^2 och y är sambandet linjärt.

```
> x<-seq(-4,4,length=50)
>
> y<-x^2+rnorm(50,sd=0.5)
> par(mfrow=c(1,2))
> plot(x,y)
> plot(x^2,y)
```



8. Som vi ser ovan kan vi skapa en mängd olika transformationer beroende på hur vi definierar $g(x)$. Vi kan tänka oss att vi skapar flera transformationer av x , vi kan kalla transformation nummer j för $z_j = g_j(x)$, alla våra transformationer blir då:

$$z_j = g_j(x) \quad j = 1, \dots, J \quad \Leftrightarrow$$

$$z_1 = g_1(x), z_2 = g_2(x), \dots, z_J = g_J(x)$$

9. När vi skapar nya variabler på detta sätt så kan vi kalla dessa för *basfunktioner* (basis functions) eller baser. Så vi skapar en eller flera basfunktioner som baseras på den ursprungliga variabeln x . Vi skriver

$$\begin{aligned} y &= f(x) + E \\ f(x) &= \beta_0 + \sum_{j=1}^J \beta_j \cdot g_j(x) \\ &= \beta_0 + \beta_1 \cdot g_1(x) + \beta_2 \cdot g_2(x) + \dots + \beta_J \cdot g_J(x) \\ &= \beta_0 + \beta_1 \cdot z_1 + \beta_2 \cdot z_2 + \dots + \beta_J \cdot z_J \end{aligned} \tag{2}$$

10. Eftersom funktionen $f(x)$ i ekvation (2) är linjär med avssende på basfunktionerna $g_j(x)$ kan vi enkelt skatta β med OLS, eller med Ridge/LASSO regression. Ekvation (2) ger oss ett ramverk för att modellera många olika typer av icke-linjära samband mellan x och y . Nedan kommer vi titta närmare på ett antal fall där vi kan använda oss av detta ramverk.

1.2 Stegfunktionen

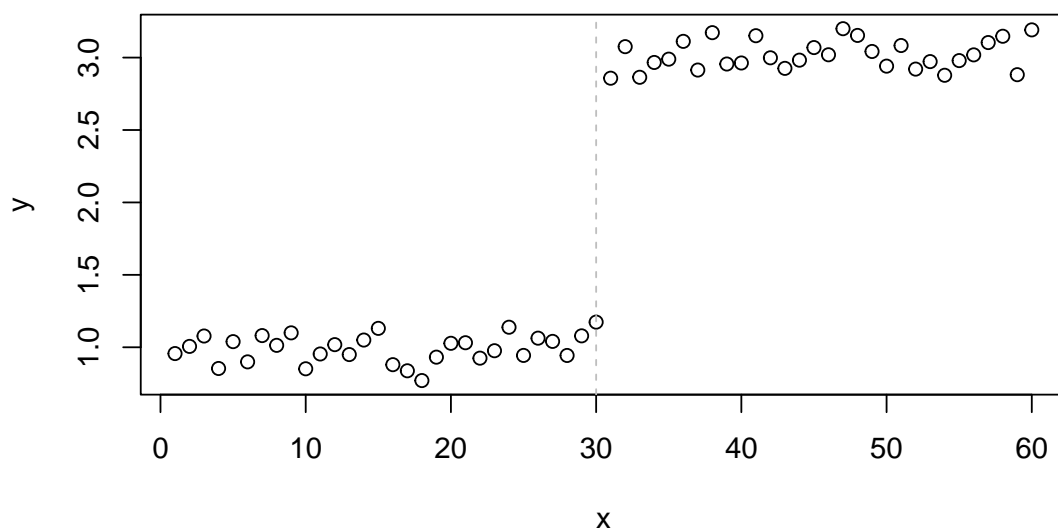
1. Stegfunktionen: En vanlig funktion inom matematiken är stegfunktionen $u(x)$, som defineras som

$$u(x) = \begin{cases} 0 & \text{om } x \leq a \\ 1 & \text{om } x > a \end{cases}$$

där a är en konstant som vi väljer. Stegfunktionen är användbar om vi har plötsliga nivåskillnader eller skiften i y .

2. Se figuren nedan. Vi har en tydlig nivåskillnad i y när $x = 30$, vilket är vid den streckade linjen.

```
> x<-1:60
> set.seed(2322)
> y<-c(rnorm(n=30,mean=1,sd=0.1),rnorm(n=30,mean=3,sd=0.1))
> plot(x,y)
> abline(v=30,lty="dashed",col="gray")
```



3. Vi testar nu att skatta en enkel linjär regression för x och y .

```
> lm_obj1<-lm(y~x)
> summary(lm_obj1)
```

Call:
lm(formula = y ~ x)

Residuals:

Min	1Q	Median	3Q	Max
-0.9331	-0.3964	-0.0457	0.3639	0.9942

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
--	----------	------------	---------	----------

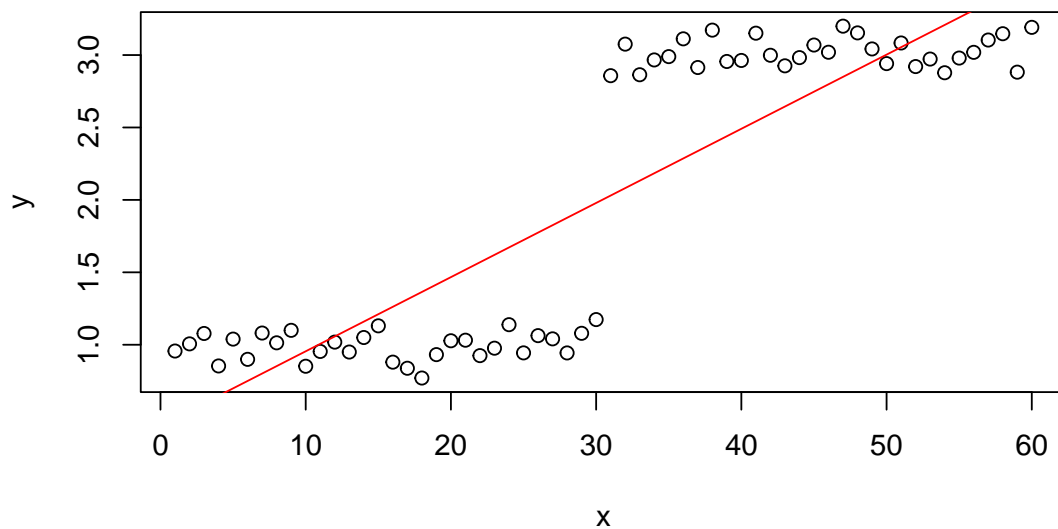
```

(Intercept) 0.44280 0.13291 3.33 0.0015 **
x           0.05120 0.00379 13.51 <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.508 on 58 degrees of freedom
Multiple R-squared: 0.759, Adjusted R-squared: 0.755
F-statistic: 183 on 1 and 58 DF, p-value: <2e-16

> y_hat<-fitted(lm_obj1)
> plot(x,y)
> lines(x,y_hat,col="red")

```



4. Vi ser att vi får signifikanta parameterar, men vi ser när vi plottar \hat{y} tillsammans med data att vi har en mycket dålig anpassning. Vi testar nu att skapa den nya variabeln

$$z = u(x) = \begin{cases} 0 & \text{om } x \leq 30 \\ 1 & \text{om } x > 30 \end{cases}$$

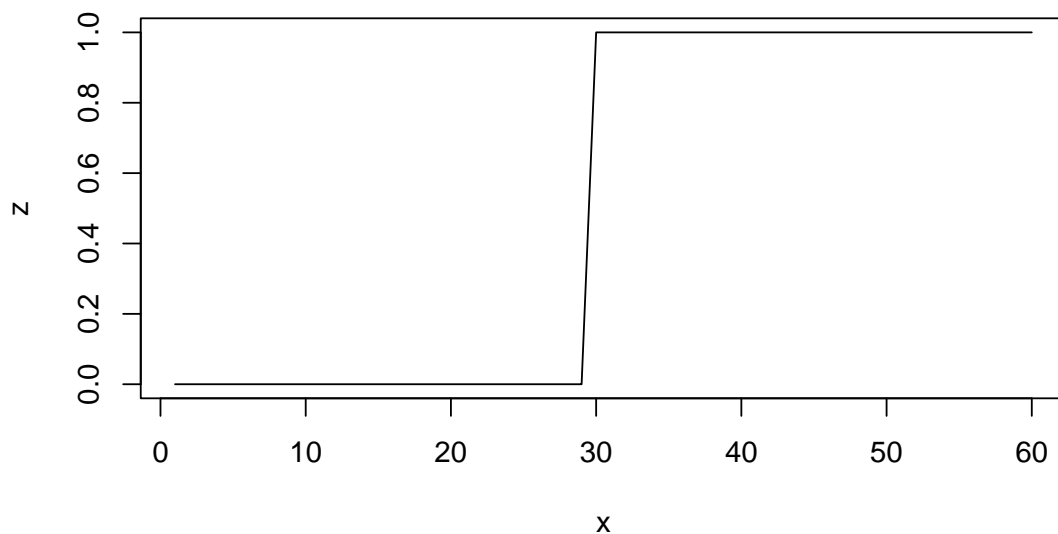
och lägger till den vår modell.

```

> z<-ifelse(x<30,0,1)
> plot(x,z,t="1",main="Stegfunktion: z = u(x)")

```

Stegfunktion: $z = u(x)$



```
> lm_obj2<-lm(y~x+z)
> summary(lm_obj2)
```

Call:

```
lm(formula = y ~ x + z)
```

Residuals:

Min	1Q	Median	3Q	Max
-1.6696	-0.0493	0.0245	0.0962	0.2675

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	0.87031	0.07271	11.97	<2e-16 ***
x	0.00765	0.00373	2.05	0.045 *
z	1.74356	0.12918	13.50	<2e-16 ***

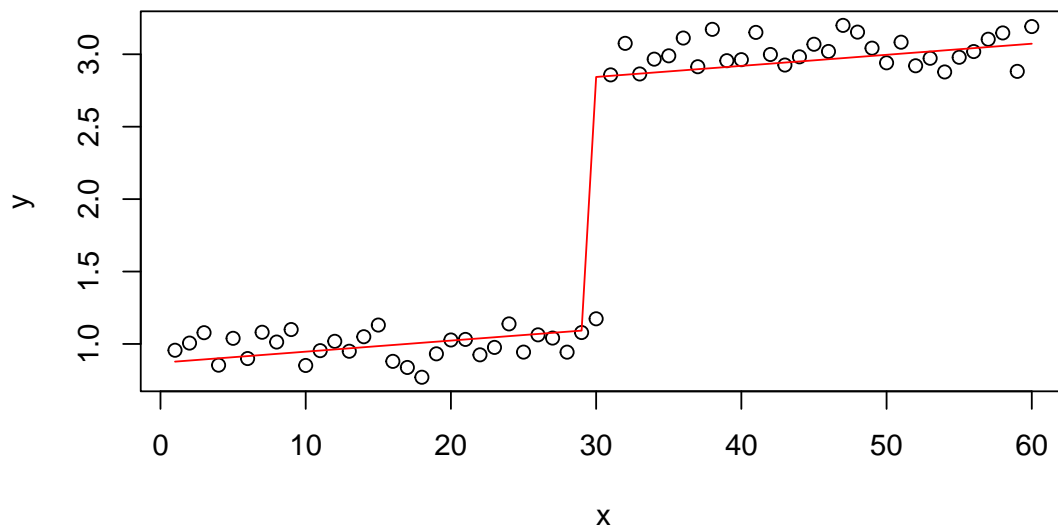
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.25 on 57 degrees of freedom

Multiple R-squared: 0.943, Adjusted R-squared: 0.941

F-statistic: 468 on 2 and 57 DF, p-value: <2e-16

```
> y_hat<-fitted(lm_obj2)
> plot(x,y)
> lines(x,y_hat,col="red")
```



5. Vi har nu en tydligt bättre anpassning för \hat{y} . Vi ser nu att x inte längre är signifikant, men att z är tydligt signifikant. Vi testar nu att bara ha med z som förklarande variabel

```
> lm_obj3<-lm(y~z)
> summary(lm_obj3)

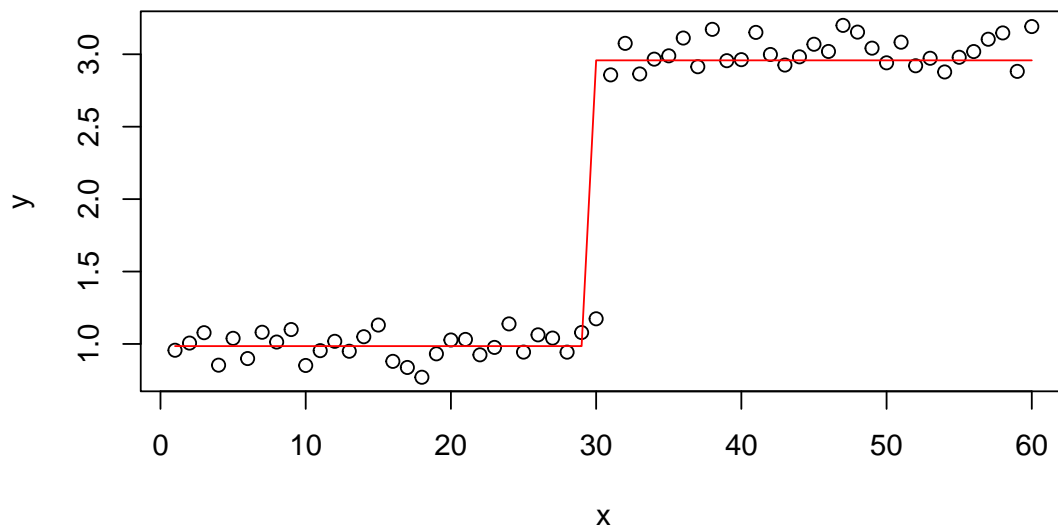
Call:
lm(formula = y ~ z)

Residuals:
    Min       1Q   Median       3Q      Max
-1.7843 -0.0414  0.0262  0.0943  0.2413

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)   0.9851     0.0478   20.6   <2e-16 ***
z             1.9731     0.0664   29.7   <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

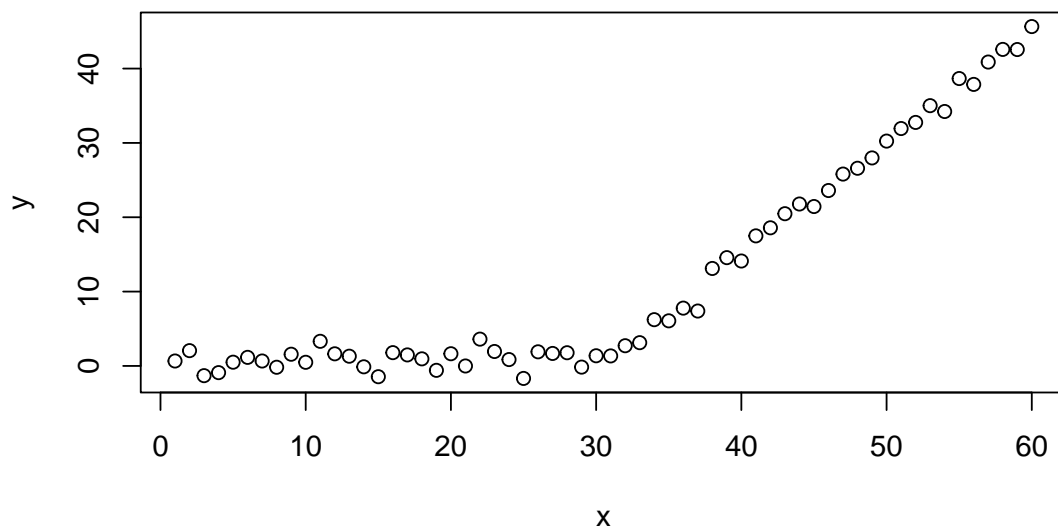
Residual standard error: 0.257 on 58 degrees of freedom
Multiple R-squared:  0.938, Adjusted R-squared:  0.937
F-statistic: 882 on 1 and 58 DF, p-value: <2e-16

> y_hat<-fitted(lm_obj3)
> plot(x,y)
> lines(x,y_hat,col="red")
```



6. Vi ser nu att vi har en ännu bättre anpassning till data, och att z är mycket signifikant. Den sista modellen som skattades är den rätta modellen, den som data generades ifrån.
7. Stegfunktioner kan hjälpa oss att modellera många olika sorters funktioner $f(x)$. Vi kan få vår modell att "bete sig" annorlunda i olika regioner av x -värden. Se datan nedan.

```
> x<-1:60
> set.seed(755)
> y<-ifelse(x<=30,1,1.5*(x[31:60]-30))+rnorm(n=60,sd=1)
> #y#<-c(rnorm(n=30,mean=1,sd=0.1),0.1*x[31:60]+rnorm(n=30,mean=0,sd=0.1))
> plot(x,y)
```



8. Här har vi en konstantfunktion fram till $x = 30$ och sen verkar vi ha en linjär funktion som med avseende på x . Vi skapar igen variabeln

$$z = u(x) = \begin{cases} 0 & \text{om } x \leq 30 \\ 1 & \text{om } x > 30 \end{cases}$$

och sen skattar vi två olika modeller:

$$y = \beta_0 + \beta_1 x + \beta_2 z + E$$

$$y = \beta_0 + \beta_1 x + \beta_2 z + \beta_3 x \cdot z + E$$

Den andra modellen har en interaktionsterm mellan x och z . Vi testat att skatta dessa modeller.

```
> x<-1:60
> z<-ifelse(x<=30,0,1)
> xz<-x*z
> model1<-lm(y~x+z)
> model2<-lm(y~x+z+xz)
>
> summary(model1)
```

Call:

```
lm(formula = y ~ x + z)
```

Residuals:

Min	1Q	Median	3Q	Max
-11.454	-4.822	-0.009	5.779	11.623

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-11.111	1.990	-5.58	6.8e-07 ***
x	0.773	0.101	7.67	2.4e-10 ***
z	-0.956	3.490	-0.27	0.79

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

Residual standard error: 6.75 on 57 degrees of freedom

Multiple R-squared: 0.795, Adjusted R-squared: 0.788

F-statistic: 111 on 2 and 57 DF, p-value: <2e-16

```
> summary(model2)
```

Call:

```
lm(formula = y ~ x + z + xz)
```

Residuals:

Min	1Q	Median	3Q	Max
-2.765	-0.874	0.235	0.699	2.591

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	0.5047	0.4371	1.15	0.25
x	0.0232	0.0246	0.94	0.35

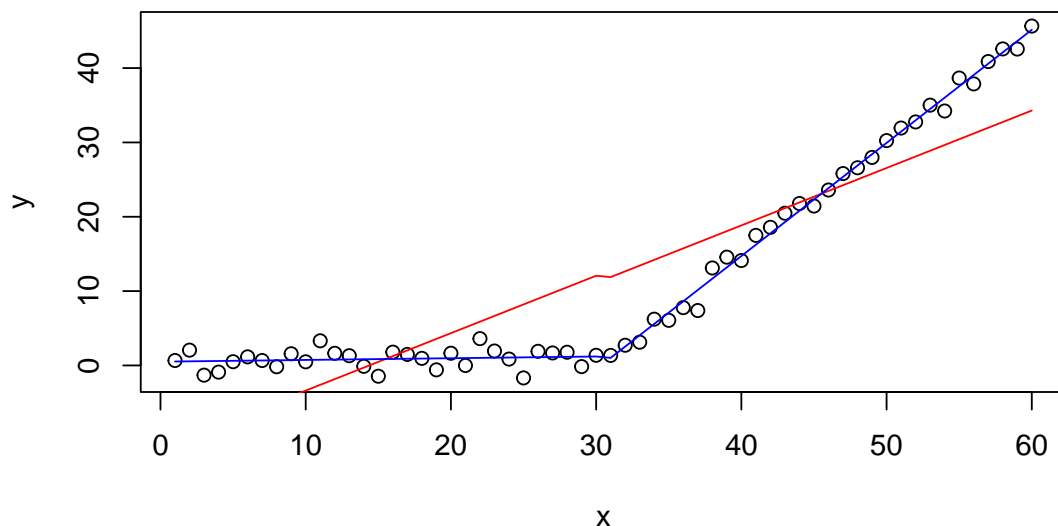
```

z          -46.6706      1.2214    -38.21    <2e-16 ***
xz           1.4988      0.0348     43.04    <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1.17 on 56 degrees of freedom
Multiple R-squared:  0.994, Adjusted R-squared:  0.994
F-statistic: 3.09e+03 on 3 and 56 DF,  p-value: <2e-16

> y_hat1<-fitted(model1)
> y_hat2<-fitted(model2)
>
> plot(x,y)
> lines(x,y_hat1,col="red") # utan interaktion
> lines(x,y_hat2,col="blue") # med interaktion

```

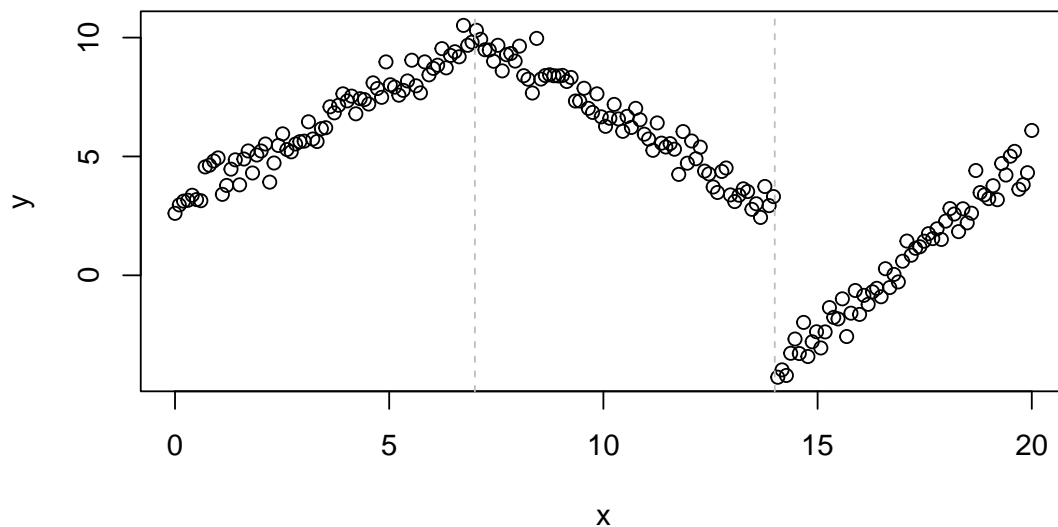


9. Den röda linjen är anpassade värden för modellen utan interaktion och den blå linjen är anpassade värden för modellen med interaktion. Här ser vi att den modellen som har med interaktionen mellan x och en lämplig stegfunktion kan modellera y väl. Vi kan ha olika lutning på linjen i olika regioner av x .
10. Vi kan generalisera detta till mer komplicerade situationer. Sen data nedan. Vi har tydliga brott vid $x = 7$ och $x = 14$. Vi kan använda den informationen när vi skapar en regressionsmodell.

```

> x<-seq(0,20,length=200)
> y1<-3+x
> y2<-17+-1*x
> y3<--25+1.5*x
> set.seed(345)
> y<-ifelse(x<7,y1,ifelse(x<14,y2,y3))+rnorm(200,sd=0.5)
> plot(x,y)
> abline(v=c(7,14),lty="dashed",col="gray")

```



11. Vi skattar en modell med två stegfunktioner, och interaktioner mellan dessa och x .

```
> # stegfunktioner:
> z1<-ifelse(x<=7,1,0)
> z2<-ifelse(x<=14,1,0)
> # interaktioner:
> xz1<-x*z1
> xz2<-x*z2
> # skatta modell:
> model<-lm(y~x+z1+z2+xz1+xz2)
> summary(model)
```

Call:

```
lm(formula = y ~ x + z1 + z2 + xz1 + xz2)
```

Residuals:

Min	1Q	Median	3Q	Max
-1.2667	-0.3531	-0.0268	0.3156	1.3956

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-25.6246	0.6463	-39.6	<2e-16 ***
x	1.5376	0.0377	40.7	<2e-16 ***
z1	-14.2725	0.3422	-41.7	<2e-16 ***
z2	42.9326	0.7214	59.5	<2e-16 ***
xz1	2.0078	0.0424	47.4	<2e-16 ***
xz2	-2.5725	0.0482	-53.4	<2e-16 ***

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

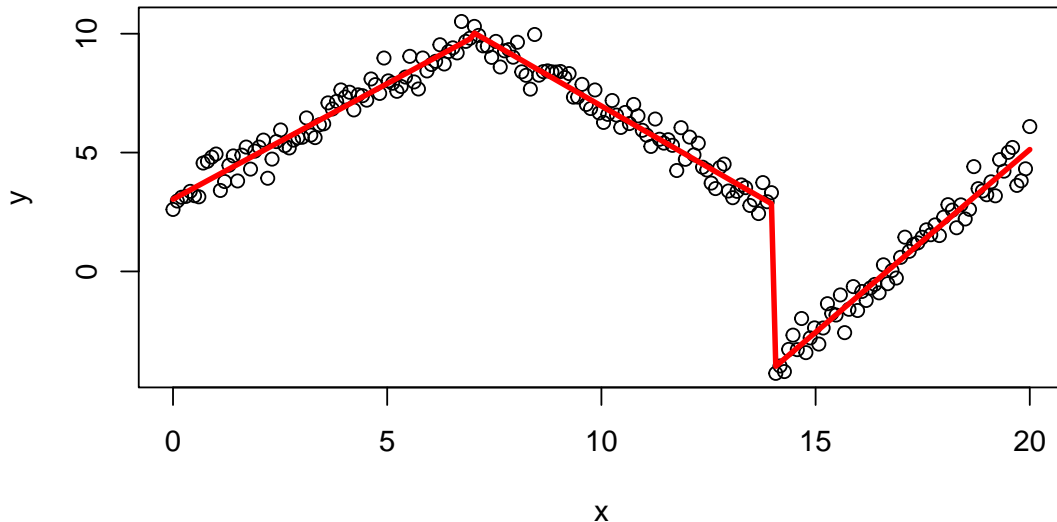
Residual standard error: 0.509 on 194 degrees of freedom

Multiple R-squared: 0.98, Adjusted R-squared: 0.979

F-statistic: 1.89e+03 on 5 and 194 DF, p-value: <2e-16

12. Vi ser att alla variabler är signifikanta och behövs för modelleringen. Vi tittar på anpassade värden nedan, och vi ser att vi har en bra anpassning av sambandet mellan x och y .

```
> plot(x,y)
> lines(x,fitted(model),col="red",lwd=3)
```



1.3 Trunkerade polynombaser

1. Trunkerade polynombaser: ovan hade vi ett fall där vi hade olika linjära samband i olika regioner av x : Ett samband för $x \leq 7$, ett för $7 < x \leq 14$ och ett för $14 < x$. Vi kan generalisera denna approach genom att använda en klass av basfunktioner som kallas trunkerade polynombaser (truncated power basis).
2. Trunkerade polynombaser är polynom, som är "avstängda" en för en del av värdena på x och sen är de bara "vanliga polynom". Man brukar ha en grupp av sådana basfunktioner. Detta är exempel på en variant av *splinemodell*. Vi kollar hur de ser ut med koden nedan. Vi använder en funktion som finns [här](#).

```
> # läs in funktion
> source(
+ "https://raw.githubusercontent.com/STIMALiU/732G12_DM/master/labs/trunc_power_basis.R"
+ )
>
> # x är en vektor med data
> # no_basis är antal basfunktioner
> basis_list<-trunc_power_basis(x = 1:120,no_basis = 6,type = "linear",show_plot = FALSE)
> str(basis_list)

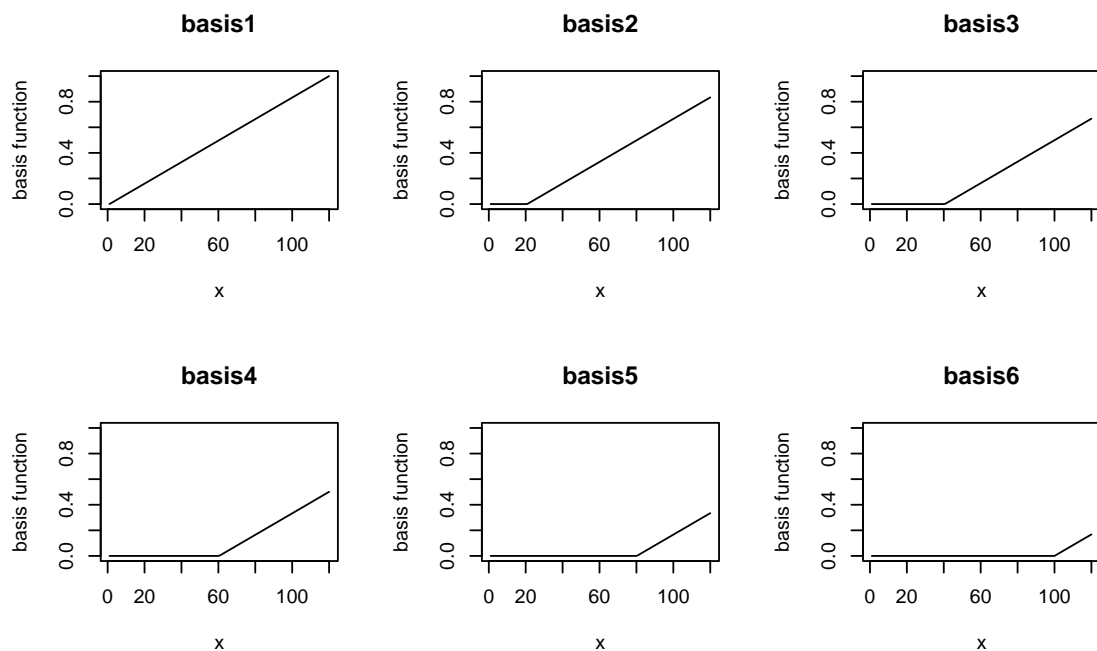
List of 2
 $ basis_mat: num [1:120, 1:6] 0 0.0084 0.0168 0.0252 0.0336 ...
  ..- attr(*, "dimnames")=List of 2
   .. ..$ : NULL
```

```

.. ..$ : chr [1:6] "basis1" "basis2" "basis3" "basis4" ...
$ x      : int [1:120] 1 2 3 4 5 6 7 8 9 10 ...

> par(mfrow=c(2,3))
> for(i in 1:6){
+   plot(basis_list$x,basis_list$basis_mat[,i],t="l",xlab="x",
+       ylab="basis function",ylim=c(0,1),main=colnames(basis_list$basis_mat)[i])
+ }

```



- Ovan ser vi att den första basfunktionen är en vanlig linjär funktion. Den andra är "avstängd" (=trunkerad) i början, sen är den en vanlig linjär funktion. Alla de andra är mer och mer trunkerade (större andel som är noll). Default är att trunkeringspunkterna är jämt spridda inom x variationsområde. Trunkeringspunkterna kallas knutar (knots), och vi kan välja dessa manuellt om vi vill. Då ger vi en numerisk vektor till argumentet `knots=`. Om vi använder dessa basfunktioner i en regressionsmodell så har vi exempel på en spline-modell. Det finns många spline-modeller, och detta är en av de enklare.
- Vi testar att använda trunkerade polynombaser för att modellera en sinusfunktion med brus.

```

> x<-1:120
> basis_list<-trunc_power_basis(x = 1:120,no_basis = 6,type = "linear",show_plot = FALSE)
> set.seed(2423)
> y<-sin(pi*x/60)+rnorm(n = 120,sd=0.1)
> plot(x,y)
> df<-data.frame(y=y,basis_list$basis_mat)
> head(df)

```

	y	basis1	basis2	basis3	basis4	basis5	basis6
1	0.017374	0.0000000	0	0	0	0	0
2	0.143701	0.0084034	0	0	0	0	0
3	0.157536	0.0168067	0	0	0	0	0
4	0.153171	0.0252101	0	0	0	0	0
5	0.223521	0.0336134	0	0	0	0	0
6	0.379642	0.0420168	0	0	0	0	0

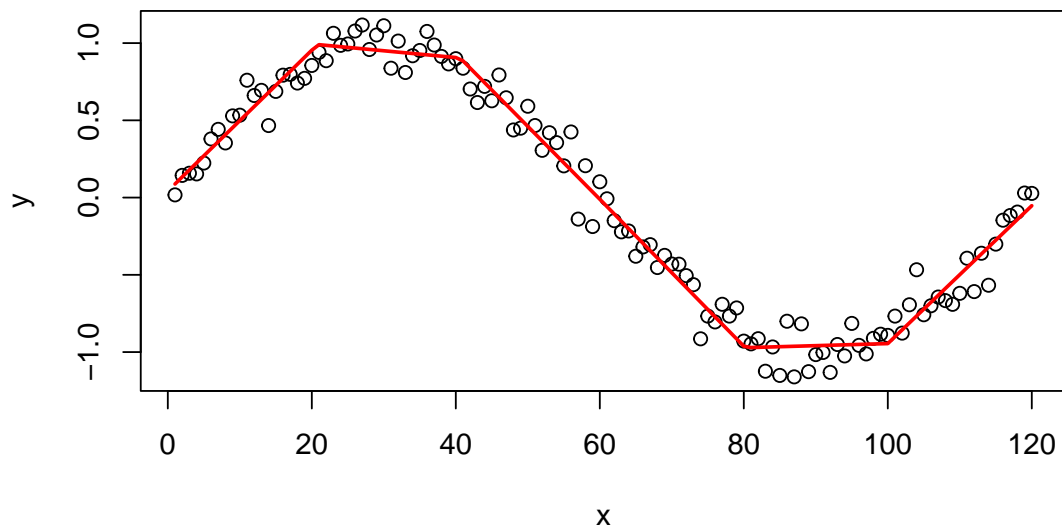
```

> tail(df)

      y basis1 basis2 basis3 basis4 basis5 basis6
115 -0.301239 0.95798 0.79160 0.62521 0.45882 0.29244 0.12605
116 -0.146450 0.96639 0.80000 0.63361 0.46723 0.30084 0.13445
117 -0.116284 0.97479 0.80840 0.64202 0.47563 0.30924 0.14286
118 -0.093748 0.98319 0.81681 0.65042 0.48403 0.31765 0.15126
119  0.029527 0.99160 0.82521 0.65882 0.49244 0.32605 0.15966
120  0.027135 1.00000 0.83361 0.66723 0.50084 0.33445 0.16807

> A<-lm(y~.,data=df)
> plot(x,y)
> lines(x,fitted(A),col="red",lwd=2)

```

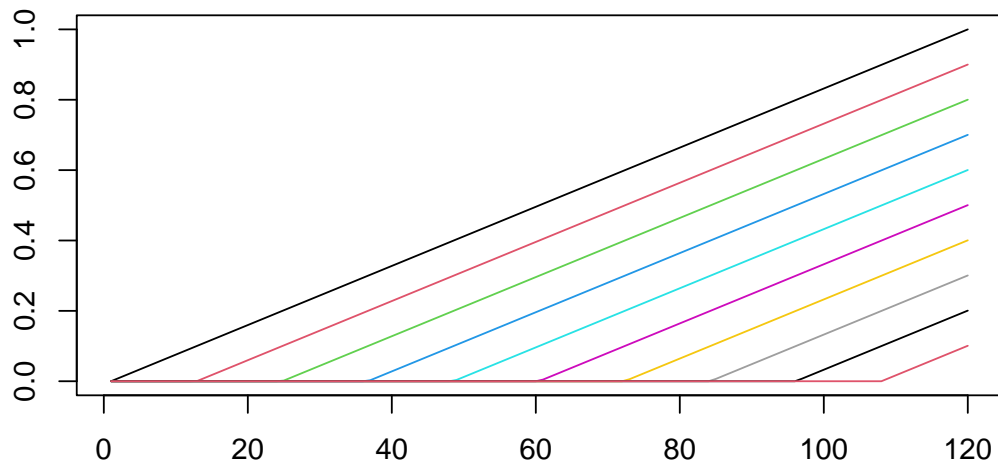


5. Vi ser att med 6 stycken linjära trunkerade polynombaser så får vi en ok anpassning av y . Ökar vi antalet basfunktioner så får vi en mer flexibel modell, men riskerar också att få överanpassning, dvs vi börjar modellera bruset i data. Testa hur anpassningen blir om vi har `no_basis = 3`, `no_basis = 10`, `no_basis = 100`. Vilken av dessa modeller är du mest nöjd med? Beräkna AIC för alla modeller, vilken har lägst värde?
6. Vi kan välja om vi vill ha linjära baser eller kvadratiske baser. De linjära kan vara bättre om det är skarpa skiftningar i y när x ändrar sig. De kvadratiske är bättre om y varierar mer mjukt när x ändrar sig.

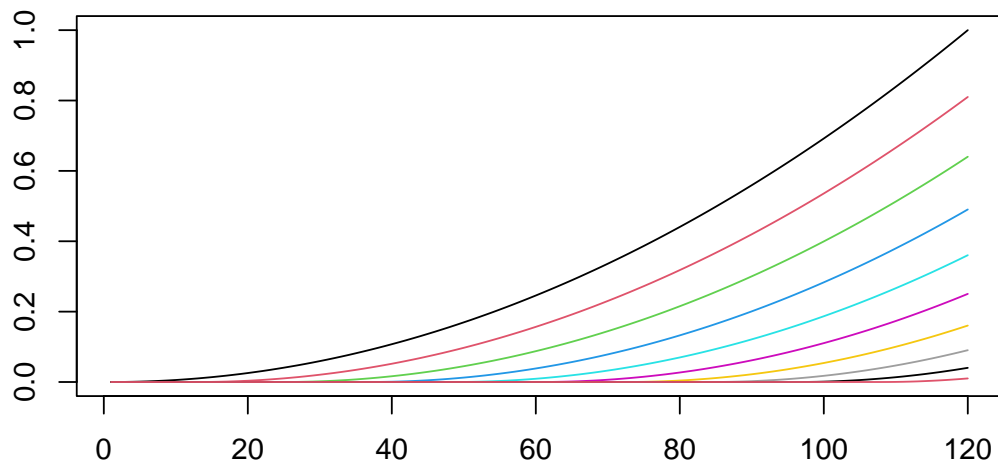
```

> # x är en vektor med data
> # no_basis är antal basfunktioner
> # vi plottar alla baserna i samma plot:
> # linjära baser:
> basis_list1<-trunc_power_basis(x = 1:120,no_basis = 10,type = "linear",show_plot = TRUE)

```



```
> # kvadratiska baser:
> basis_list2<-trunc_power_basis(x = 1:120,no_basis = 10,type = "quadratic",show_plot = TRUE)
```



7. Testa att modellera sinusfunktionen ovan med kvadratiska trunkerade polynombaser. Testa låta `no_basis` vara 5, 10 och 30. Vilken anpassning är ni mest nöjda med?
8. Har vi många basfunktioner i vår "fulla modell" så är det ofta bra att använda metoder för automatiskt modellval. Här kan krympningsmetoder som ridge eller LASSO användas. Andra alternativ är framåtvalsmetoden eller Stepwise forward selection³. Har vi många basfunktioner och/eller många observationer så rekommenderas paketet `bigstep`, se [här](#).

³Varför är inte backward elimination så bra här?

1.4 Consinusbaser

En annan variant av basfunktioner som kan modellera icke-linjära funktioner är cosinusbaser. De definieras som

$$\begin{aligned}x_{min} &= \min(x) & x_{max} &= \max(x) & L &= x_{max} - x_{min} \\ g_h(x) &= \cos\left(\frac{\pi \cdot h \cdot (x - x_{min})}{L}\right) & h &= 1, 2, \dots, H\end{aligned}\quad (3)$$

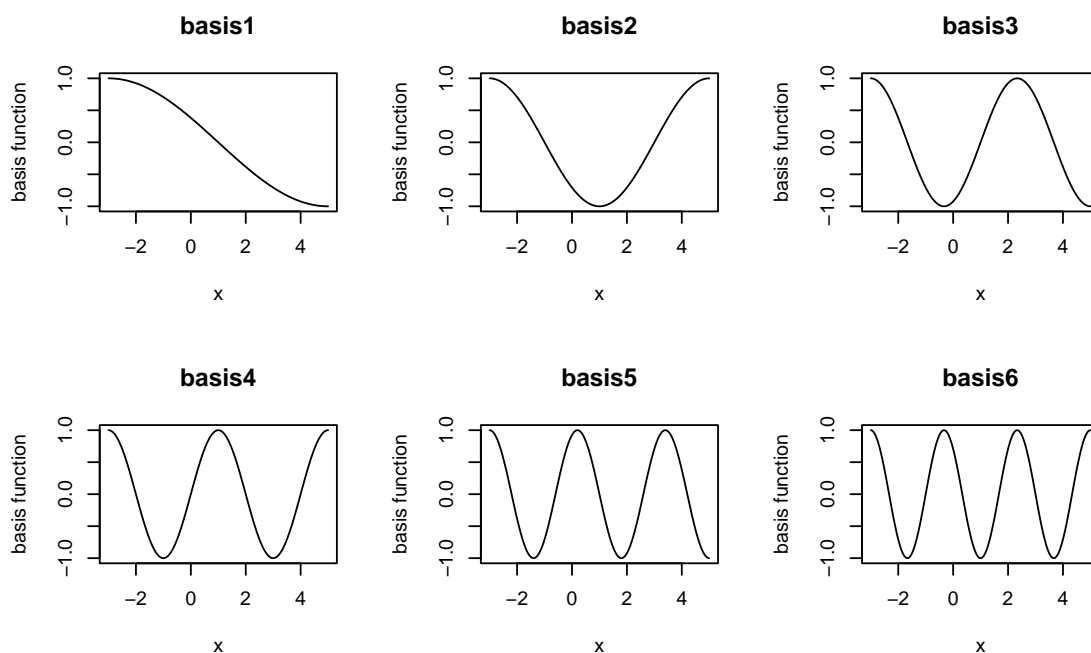
där x är en vektor med vår förklarande variabel, x_{min} är det minsta värdet av x och x_{max} är det största värdet av x . h är ordningen för basen, och H är den högsta ordningen. Vi kan likna h med graden i ett polynom. Så vi kan givet vår förklarande variabel skapa H stycken basfunktioner baserat på ekvation (3). Desto högre H är desto mer flexibel och icke-linjär funktion kan vi modellera med dessa baser. Samtidigt ökar risken för överanpassning ju högre H som vi har. En fördel med dessa basfunktioner är att de är okorrelerade med varandra.

1. Vi testar att skapa några cosinusbaser och ser hur de ser ut. Vi använder en funktion som finns [här](#).

```
> source(
+ "https://raw.githubusercontent.com/STIMALiU/732G12_DM/master/labs/cosine_basis.R"
+ )
>
> x<-seq(-3,5,length=300)
> # order: högsta ordningen på cosinusbaserna
> cos_list<-cosine_basis(x = x,order = 6)
> dim(cos_list$basis_mat)

[1] 300    6

> par(mfrow=c(2,3))
> for(i in 1:6){
+   plot(cos_list$x,cos_list$basis_mat[,i],t="l",xlab="x",ylab="basis function",ylim=c(-1,1),
+ }
```



2. Vi ser att den första basen har en halv period mellan x_{min} och x_{max} . Sen ju högre ordning som basen har desto fler svängningar har den mellan x_{min} och x_{max} . Vi testar att använda dessa basfunktioner

för att modellera en observerad y -variabel.

```
> n_obs<-300
> x<-1:n_obs
> set.seed(4325)
> y<-arima.sim(model = list(ar=0.999),n = n_obs,sd=1)
> plot(x,y,t="p")
>
> # skapar baser:
> cos_list<-cosine_basis(x = x,order = 6)
> # anpassar en modell
> df<-data.frame(y=y,cos_list$basis_mat)
> A<-lm(y~.,data = df)
> summary(A)
```

Call:

lm(formula = y ~ ., data = df)

Residuals:

Min	1Q	Median	3Q	Max
-4.800	-1.288	-0.129	1.322	4.958

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	34.710	0.111	312.54	< 2e-16 ***
basis1	-4.135	0.157	-26.37	< 2e-16 ***
basis2	-2.106	0.157	-13.43	< 2e-16 ***
basis3	-1.323	0.157	-8.44	1.5e-15 ***
basis4	-1.368	0.157	-8.72	< 2e-16 ***
basis5	-1.081	0.157	-6.90	3.3e-11 ***
basis6	-1.064	0.157	-6.79	6.4e-11 ***

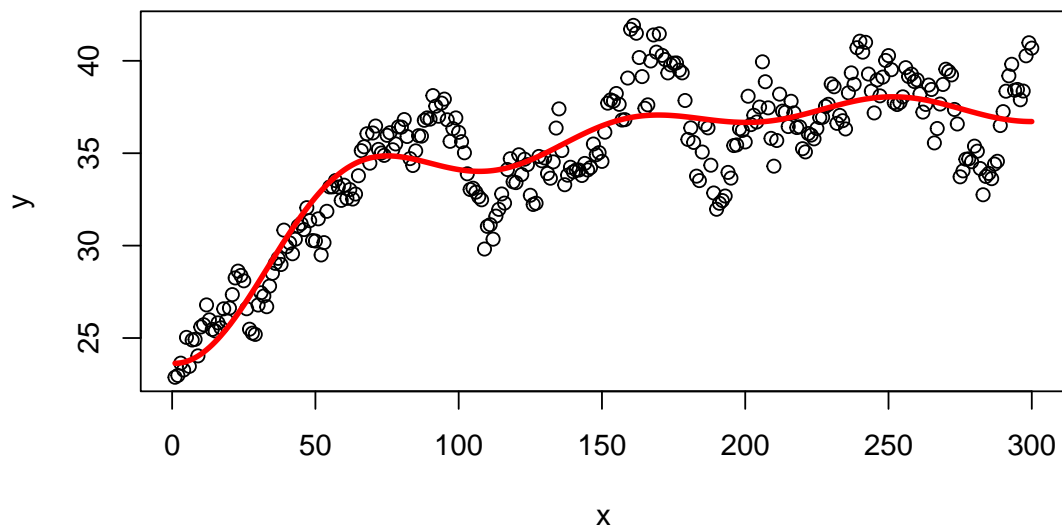
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1.92 on 293 degrees of freedom

Multiple R-squared: 0.794, Adjusted R-squared: 0.789

F-statistic: 188 on 6 and 293 DF, p-value: <2e-16

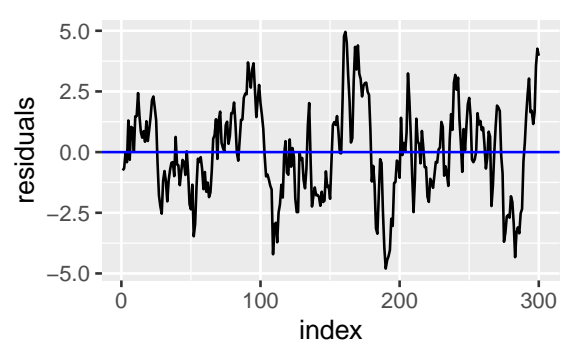
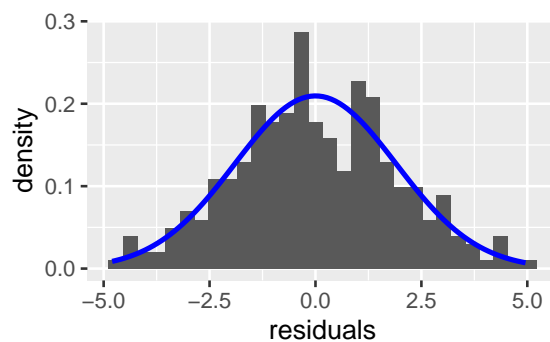
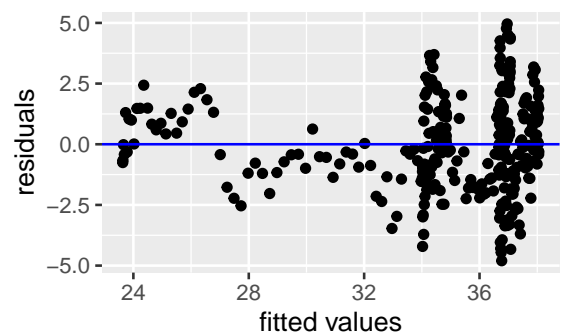
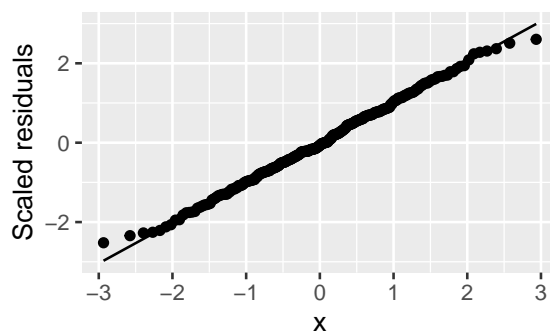
```
> y_hat<-fitted(A)
> plot(x,y,t="p")
> lines(x,y_hat,col="red",lwd=3)
```



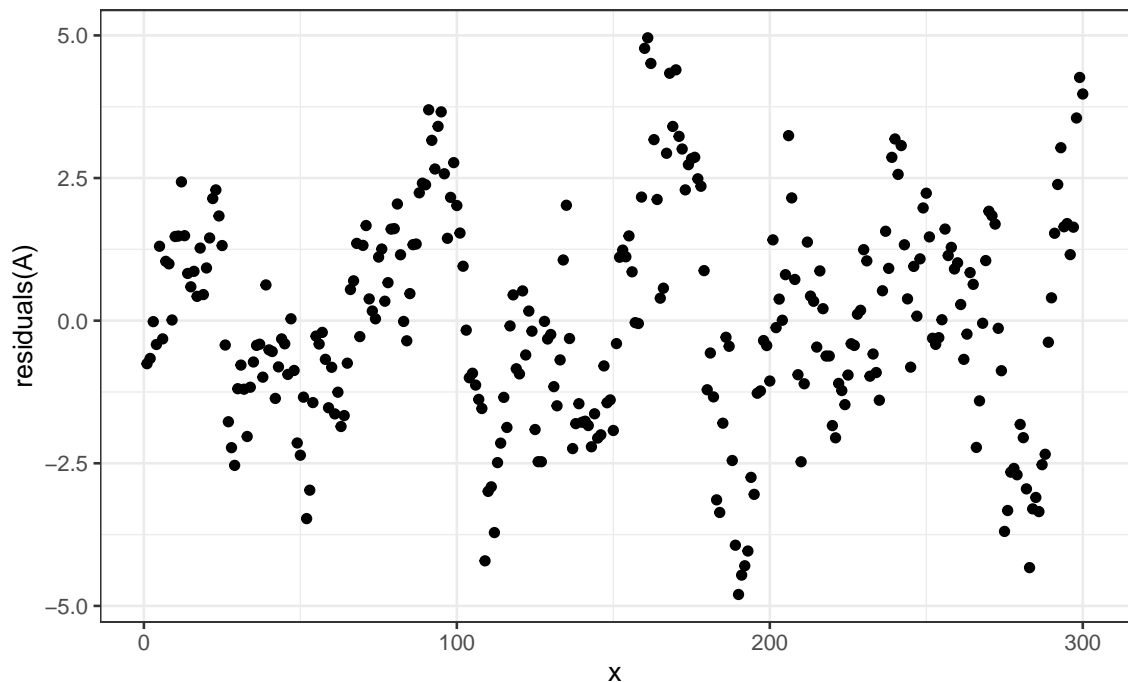
3. Vi ser att vi har en bättre anpassning än en vanlig linjär modell med bara x som kovariat. Dock så kan man önska en bättre anpassning. Vi gör residualanalys:

```
> suppressWarnings({
+   lm_diagnostics(A)
+ })
```

`stat_bin()` using `bins = 30`. Pick better value with `binwidth`.



```
> # förklarande variabeln mot residualer
> qqplot(x,residuals(A))+theme_bw()
```



4. Vi ser att det finns beroende kvar mellan x och y . Tex ser vi tydliga mönster i fitted values mot residuals.
5. Anpassa nu en linjär regression med cosinusbaser med `order = 10` och `order = 30` till data ovan. Hur blir anpassningen?

1.5 ISL

1. Gå igenom Lab 7.8.1 i ISL
2. Vad är skillnaden mellan följande tre funktioner för polynomregression?

```
> y~poly(x,4)
> y~poly(x,4,raw=TRUE)
> y~x+I(x^2)+I(x^3)+I(x^4)
```

1.6 Sammanfattning

1. När vi ska modellera icke-linjära samband mellan x och y , så är det ofta bra att skapa transformationer av x , olika alternativ:
 - (a) Ofta vill vi använda en metod för att "automatiskt" transformera x , exempel är splines eller lokal regression → **detta är ofta att föredra!**
 - (b) Vi kan välja ett antal manuella transformationer av x som passar vårt problem, tex $\log(x)$, $\sin(\pi \cdot x)$, x^2 , interaktioner, stegfunktioner etc.
 - i. Testa att plotta den nya transformerade variabeln mot y . Om du ser ett linjärt mönster i plotten så är troligen transformationen lämplig. Detta funkar ofta bra om vi har "ett tydligt mönster" mellan x och y .
2. När vi använder oss av många basfunktioner, som också kan vara komplicerade, då är det bra att tänka på:

- (a) Vi är sällan intresserade av inferens för enskilda basfunktioner, utan ofta fokuserar vi mer på vilken *funktion* som basfunktionerna tillsammans skapar. Vi fokuserar på hur bra de anpassade värden blir för modellen.
- (b) Det finns ju fler och ju mer avancerade basfunktioner som vi använder, desto större är risken för *överanpassning*. Så ofta vill vi dela in data i en träningsdel och en valideringsdel, och sen beräkna MSE⁴ för båda dessa delar. Då kan vi jämföra olika mängder av basfunktioner, och se vilken som har lägst MSE på valideringsdata.
- (c) Det är bra att plotta anpassade värden för att se hur bra de modellerar y . Det är viktigt att göra residualanalys för att se om vi har lyckats modellera relationen mellan x och y på ett bra sätt. Det är bra att göra residualanalys på både träningsdelen och valideringsdelen.
- (d) Vi kan använda olika metoder för automatiskt modellval.
 - i. Ofta är krympningsmetoder att föredra (Ridge, LASSO etc)
 - ii. Andra alternativ: framåtvalsmetoden eller Stepwise forward selection. Har vi många basfunktioner och/eller många observationer så rekommenderas paketet **bigstep**, se [här](#).

2 Splines och Lokal regression

1. Gör Lab 7.8.2 i ISL
2. Ladda in datamaterialet `lab2_data_1.csv`. Välj knutarna i 1, 2, 3, och 4. Testa nu att skatta vanliga splines (`bs()`) med `degree` 1, 2, 3 och 4. Vad blir skillnaden i resultat? Vilken ser bäst ut?
3. Använd samma datamaterial och ändra till natural splines (`ns()`). Vad blir skillnaden? Lägg till knutar i 0.5 och 4.5 och jämför.
4. Testa andra platser för knutarna och använd korsvalidering för att hitta den bästa modellen. Här finns en R-fil som visar hur man kan använda korsvalidering för olika modeller: [länk](#)

3 Del 3: GAM

1. Gör Lab 7.8.3 i ISL
2. Här finns mycket material om GAMs: [GAMs in R by Noam Ross](#), detta är en fri webbkurs om GAMs med R-paketet `mgcv`.
 - Rekommenderas att gå igenom delarna "1 - Introduction to Generalized Additive Models" och "2 - Interpreting and Visualizing GAMs"
3. Kombinera GAM och variabel selektion, det finns olika paket som kan göra olika varianter av variabelselektion när vi arbetar med GAMs:
 - (a) `gamsel`, här finns en vingette: [länk](#)
 - i. `gamsel` kan automatiskt välja om en förklarande variabel ska vara tas bort, ha ett linjärt samband eller ett icke-linjärt samband (=splines) med y .
 - (b) `step.gam()` funktion i paketet `gam`
 - (c) `gam.selection()` funktion i paketet `mgcv`

4 Del 4: Email Spam

1. Ladda in `spambase.csv` datasetet och bekanta dig med det. Vi vill skapa en modell som predikterar **spam** (0 eller 1) givet de förklarande variablerna. Då våra förklarande är tungsvansade kan en log-transformation fungera bra ($\log(x + 0.1)$) där vi lägger på 0.1 för att undvika $\log(0)$ problem.
2. Dela upp datamaterialet i 70% träningsdata och 30% valideringsdata. Undersök så att klassbalansen i är liknande i träningsdata och valideringsdata.

⁴Eller annat lämpligt utvärderingsmått.

3. Börja med att anpassa en vanlig multipel logistisk regression till datamaterialet. Vilket klassificeringsfel får du på tränings- och valideringsdata?
4. Anpassa nu en GAM modell med kubiska natural splines med 5 frihetsgrader för varje förklarande variabel. Vilka klassificeringsfel får du nu?
5. Testa med andra ordningar av frihetsgrader för dina splines och se om du kan få till en bättre modell.