

SAMBANDSANALYS

REGRESSION OCH KORRELATION
ORIENTERING OM TIDSSERIER

HT 2012



LUNDS UNIVERSITET

Lunds Tekniska Högskola

Matematikcentrum
Matematisk statistik

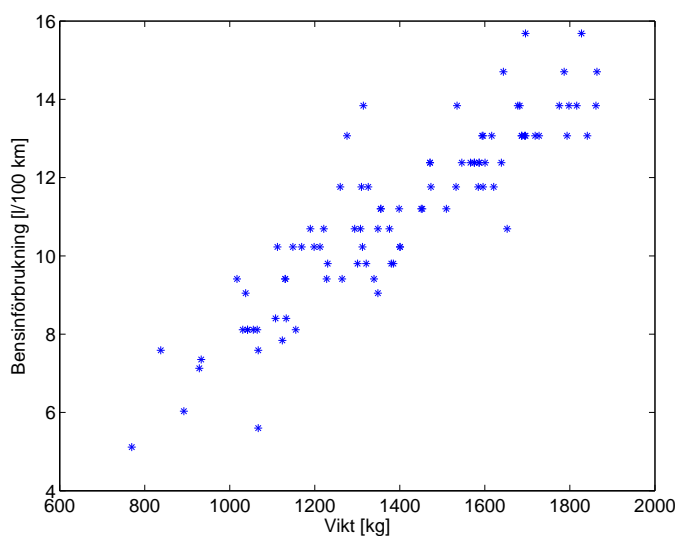
Innehåll

1	Samband mellan två eller flera variabler	3
2	Enkel linjär regression	5
2.1	Intressanta frågeställningar	5
2.2	Modellantaganden	6
2.3	Skattningar av parametrarna α , β och σ	7
2.4	Konfidensintervall för α och β	8
2.5	Skattning av punkt på linjen	9
2.6	Prediktionsintervall för observationer	10
2.7	Kalibreringsintervall	12
2.8	Modellvalidering	12
2.8.1	Residualanalys	12
2.8.2	Är β signifikant?	15
2.9	Förklaringsgrad	15
2.10	Outliers	16
2.11	Linjärisering av några icke linjära samband	16
2.12	Jämförelse av två lutningar	17
3	Multipel linjär regression på matrisform	20
4	Korrelationsanalys	22
4.1	Mått på samband	22
4.2	Test av samband	23
4.3	Var försiktig med korrelationskoefficienten!	24
4.4	Anknytning till linjär regression	25
5	Tidsserier	26
5.1	Syftet med analysen	26
5.2	Beskrivning av tidsserien	27
5.2.1	Komponentuppdelning	27
5.2.2	Beroende i tidsserien	28
5.2.3	Skattning av autokorrelationsfunktionen	29
5.2.4	Matlabkommandon för skattning av autokorrelationsfunktionen	30
5.3	Modeller	32
5.3.1	AR(1)-processer	32
5.3.2	Simulering av AR(1)-processer i Matlab	34
5.4	Beroende mätningar påverkar analysen	34
5.4.1	Beroende data påverkar trendanalysen!	35
5.5	Läsa mer om trendanalys och tidsserier	37
6	Mer om trendanalys	38
6.1	Mann-Kendalls test	38
6.2	Skattning av trenden	39
6.3	Seasonal Kendall test	40
7	Appendix: ML- och MK skattningar av parametrarna i enkel linjär regression	43
7.1	Några hjälpresultat	43
7.2	Punktskattningar	43
7.3	Skattningarnas fördelning	44

1 Samband mellan två eller flera variabler

Det är ganska vanligt att man gör mätningar på två eller flera variabler och vill undersöka om det finns något samband mellan dem. Vi presenterar två exempel:

Exempel 1.1. För ett slumpmässigt urval av bilar noterar man y -bensinförbrukning i stadskörning (l/100 km) och x -vikt (kg). Data beskrivs i figur 1 där y plottats mot x . \square



Figur 1: Ett slumpmässigt urval av bilar där $y =$ "bensinförbrukning i stadskörning" är plottad mot $x =$ "vikt".

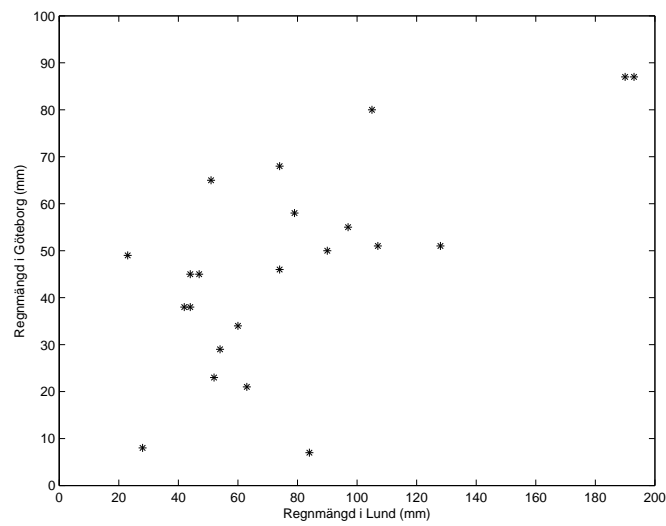
I detta exempel är det rimligt att tänka sig att y -bensinförbrukning påverkas av x -vikt (och inte tvärt om!). Vi kan alltså försöka beskriva y som en funktion av x , analysen måste naturligtvis också ta hänsyn till att mätningarna påverkas av en slumpmässig störning. Vi gör en *regressionsanalys* där y är responsvariabeln medan x är den förklarande variabeln. Ibland kallas även y för den beroende variabeln medan x är den oberoende variabeln:

$$\underbrace{y}_{\text{responsvariabel}} = \underbrace{f(x)}_{\text{regressionsfunktion med förklarande variabel } x} + \underbrace{\text{"slump"}}_{\text{s.v. med fördelning}}$$

När regressionsfunktionen $f(x)$ är linjär med avseende på sina parametrar har vi linjär regression. Från figuren verkar det rimligt att tänka sig ett linjärt samband mellan x och y som beskriver hur stor bensinförbrukning en "medelbil" av en viss vikt har. Om man, som i vårt exempel, har enbart en förklarande variabel, x , talar man om *enkel linjär regression*. Hela nästa avsnitt kommer att behandla denna viktiga situation.

Exempel 1.2. Månadsnederbörden, d.v.s. den totala mängden nederbörd (mm) under en månad, noterades i Göteborg och Lund under åren 2005 och 2006. I figur 2 markerar varje punkt en månad där Göteborgs nederbörd avläses på y -axeln och Lunds på x -axeln. \square

Här är det inte självklart att någon av de två uppmätta variablerna kan beskrivas som en funktion av den andra. Variablerna är "likvärdiga" eftersom vi lika gärna skulle kunna byta variabel på axlarna och placera Lundamätningarna på y -axeln och Göteborgsmätningarna på x -axeln. I denna situation är det olämpligt att använda regression, man får nöja sig med att beskriva graden av samband i en *korrelationsanalys*. Vi kommer att studera detta närmare i avsnitt 4.



Figur 2: Månadsvisa mätningar av nederbörden (mm) där $y =$ ”nederbörd i Göteborg” är plottad mot $x =$ ”nederbörd i Lund”.

2 Enkel linjär regression

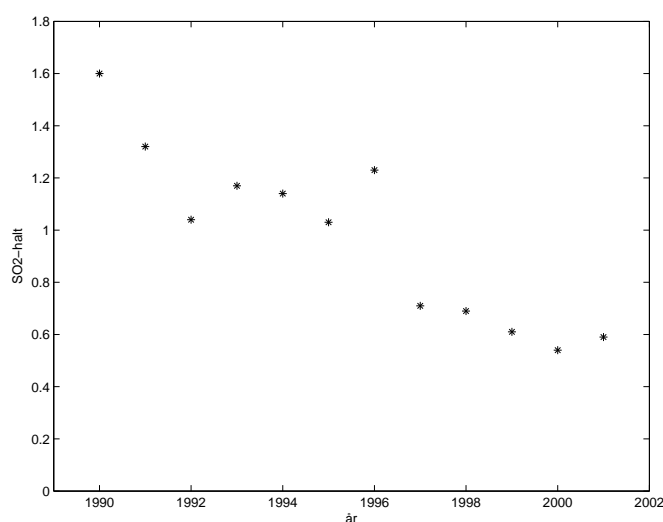
I enkel linjär regression studerar vi en variabel y som beror linjärt av en variabel x men samtidigt har en slumpmässig störning eller avvikelse:

$$Y_i = \alpha + \beta x_i + \varepsilon_i,$$

där ε_i är den slumpmässiga avvikelsen från linjen.

I detta avsnitt kommer vi illustrera teorin med hjälp av två dataset: mätningarna från exempel 1.1 om bensinförbrukning hos bilar samt mätningar av SO_2 -halt i luft.

Exempel 2.1. Inom miljöövervakningsprogrammet EMEP har man under en lång period mätt årsmedelhalter av SO_2 ($\mu g/m^3$) i Hoburgen på Gotland. I figur 3 visas halterna under åren 1990-2001 (källa: IVL Svenska Miljöinstitutet AB, www.ivl.se). □



Figur 3: Mätningar vid Hoburgen på Gotland $y = "SO_2\text{-halt}"$ ($\mu g/m^3$) är plottad mot $x = "år"$.

2.1 Intressanta frågeställningar

Det finns en mängd frågeställningar kring den beskrivna situationen som är intressanta:

- Hur ska vi skatta α och β i regressionslinjen $y = \alpha + \beta x$? Lutningen β beskriver hur mycket y ändras då x ökar med en enhet: hur mycket ökar bensinförbrukningen då vikten hos en bil ökar med ett kg? Speciellt intressant är det att undersöka om $\beta = 0$ eftersom det innebär att regressions sambandet då kan reduceras till $y = \alpha$, d.v.s. att y inte beror av x . I data från Hoburgen innebär ett $\beta \neq 0$ att det finns en trend i SO_2 -halt.
- Hur stor är variationen kring linjen? Eftersom ε_i beskriver den slumpmässiga avvikelsen från linjen motsvarar det att undersöka hur stor denna avvikelse tenderar att vara - ett mått på detta är $D(\varepsilon_i)$ som vi betecknar σ .
- Givet ett x_0 , vad är det **förväntade** värdet på Y ? Vi söker alltså $\mu_0 = \alpha + \beta \cdot x_0$, linjens läge i punkten x_0 . I bil exemplet kan vi t.ex. vara intresserade av hur stor bensinförbrukningen är i genomsnitt hos bilar som väger 1200 kg. I Hoburgsdata vad förväntad SO_2 -halt var 1994.
- Skilj den föregående frågeställningen från följande: Givet ett x_0 , vad är en **enstaka** observation av Y , Y_0 ? Vi vill göra en prediktion av Y -värdet. Det kan t.ex. gälla en prognos av Y för något framtida

värde på x . Om vi har en bil som väger 1200 kg, är vi nu intresserade av hur stor bensinförbrukningen är för detta exemplar. I SO_2 -exemplet kan vi vilja prediktera halten för år 2002 - inom vilket intervall är det troligt att kommer den att hamna?

- Hur bra passar modellen till data? Är det lämpligt att beskriva sambandet med en linjär funktion eller borde vi ansätta något annat? Denna frågeställning bör man studera först - det är naturligtvis viktigt att den antagna modellen stämmer någorlunda till data innan man detaljstuderar den.
- Hur mycket av den totala variationen i y -led har vi förklarat med modellen? Man kan inte räkna med att modellen ska förklara all variation som finns i mätningarna. Bensinförbrukningen hos en bil beror inte enbart på bilens vikt utan påverkas - förutom av slumpmässig variation - av en mängd andra variabler. Hur stor andel av variation i bensinförbrukning kan beskrivas med hjälp av bilars vikt och hur stor andel av variationen återstår att beskriva? Den återstående variationen kanske delvis kan förklaras med hjälp av andra variabler?

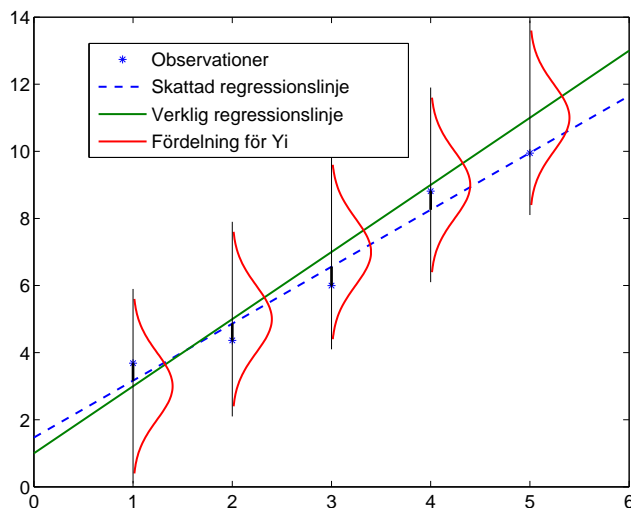
För att kunna hantera dessa frågor gör vi vissa antaganden om den linjära modellen och om våra mätningar $(x_1, y_1), \dots, (x_n, y_n)$.

2.2 Modellantaganden

Vi använder följande modell där y_i är n st oberoende observationer av

$$Y_i = \alpha + \beta x_i + \varepsilon_i, \quad \text{där } \varepsilon_i \in N(0, \sigma), \text{ oberoende av varandra}$$

så observationerna är $Y_i \in N(\alpha + \beta x_i, \sigma) = N(\mu_i, \sigma)$, dvs de är normalfördelade med väntevärde på den okända regressionslinjen $\mu(x) = \alpha + \beta x$ och med samma standardavvikelse σ som avvikelserna ε_i kring linjen har; se figur 4.



Figur 4: Sann regressionslinje, observationer och skattad regressionslinje. Residualerna är markerade som de lodräta avstånden mellan observationerna och den skattade regressionslinjen.

Modellen ovan är beskriven i "kortform", några förklaringar och kommentarer till den:

- Vi tänker oss att x -värdena är fixa eller uppmätta med ett försumbart mätfel - ofta kan vi själva välja vilka x -värden vi vill studera. Den slumpmässiga variation vi vill modellera finns enbart i y -led. I bil-exemplet anses vikten hos en bil inte ha någon större variation; likaså är det uppenbart att x -variabeln i Hoburgsexemplet - årtalen - är fixa.

- Tidigare har vi haft modeller där mätningarna är observationer av stokastiska variabler ξ_i , vilka hade samma väntevärde μ , men nu är observationernas väntevärde en linjär funktion av x . Beteckningen Y_i är också en naturligare beteckning för den stokastiska variabeln.
- Att de slumpmässiga avvikelserna från linjen, $\varepsilon_1, \dots, \varepsilon_n$ är oberoende innebär t.ex. att om en avvikelse råkar bli stor (liten) vid ett visst x -värde ska det inte påverka hur avvikelsen blir vid något annat x -värde. Om SO_2 -halten år 1991 är lägre än vad som förväntades enligt linjens läge vid denna tidpunkt ska detta alltså inte påverka hur halten avviker från linjens läge vid t.ex. år 1992.
- För ett fixt x -värde kommer motsvarande y -mätningar att vara normalfördelade kring linjen och standardavvikelsen i den fördelningen är σ ; se figur 4. Om vi t.ex. slumpmässigt väljer ut ett antal bilar som alla har vikt 1400 kg och mäter deras bensinförbrukning kommer förbrukningen att fördela sig enligt en normalfördelning med väntevärde $\alpha + \beta \cdot 1400$ och standardavvikelse σ .
- Observera att vi tänker oss att spridningen i normalfördelningarna är den samma oavsett värde på x , d.v.s. σ är konstant. Det innebär t.ex. att modellen inte tillåter att spridningen kring linjen ändrar sig då x -värdet ändras. Det är inte ovanligt i många sammanhang att y -mätningarna uppvisar en större spridning med ökande värde på x ; för denna situation kan vi alltså inte direkt använda oss av ovanstående modell.

2.3 Skattningar av parametrarna α, β och σ

För att skatta parametrarna α och β används minsta kvadrat-metoden (MK-metoden). Skattningarna och deras fördelning härleds i appendix i avsnitt 7, här presenteras enbart resultaten.

MK-skattningarna av regressionslinjens lutning, β , och intercept, α , ges av

$$\beta^* = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{S_{xy}}{S_{xx}}, \quad \alpha^* = \bar{y} - \beta^* \bar{x}.$$

Eftersom β^* är en linjär funktion av observationerna Y_i ($\beta^* = \sum c_i Y_i$ där $c_i = (x_i - \bar{x})/S_{xx}$), och även α^* en linjär funktion av β^* och observationerna, är dessa skattningar normalfördelade med väntevärde och standardavvikelse enligt

$$\beta^* \in N\left(\beta, \frac{\sigma}{\sqrt{S_{xx}}}\right), \quad \alpha^* \in N\left(\alpha, \sigma \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}}\right).$$

De två skattningarna är dock *inte* oberoende av varandra. Man kan däremot visa att β^* och \bar{Y} är oberoende¹ av varandra.

Då man ska skatta variansen σ^2 visar det sig lämpligt att studera modellens s.k. *residualer*, r_1, \dots, r_n där

$$r_i = y_i - (\alpha^* + \beta^* x_i), \quad i = 1, \dots, n,$$

är residualen för x_i och motsvarar den lodräta avvikelsen mellan det observerade värdet y_i och den skattade linjen, se figur 4. Residualen r_i är ett närmevärde till den slumpmässiga avvikelsen ε_i och eftersom σ^2 är ett mått på spridningen hos ε_i är det rimligt att residualerna kan användas när vi vill skatta variansen.

En väntevärdesriktig skattning av variansen ges av

$$(\sigma^2)^* = s^2 = \frac{Q_0}{n-2}$$

där Q_0 är residualkvadratsumman

$$Q_0 = \sum_{i=1}^n (y_i - \alpha^* - \beta^* x_i)^2 = \sum_{i=1}^n r_i^2 = S_{yy} - \frac{S_{xy}^2}{S_{xx}}.$$

¹Vi visar inte här att β^* och \bar{Y} är oberoende av varandra, men det faktum att regressionslinjen alltid går genom punkten (\bar{x}, \bar{y}) gör det kanske troligt; om β över- eller underskattas påverkas inte \bar{Y} av detta.

För att räkna ut kvadratsummorna S_{xx} , S_{yy} och S_{xy} ”för hand” kan man ha användning av sambanden

$$S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2$$

$$S_{yy} = \sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n y_i^2 - \frac{1}{n} \left(\sum_{i=1}^n y_i \right)^2$$

$$S_{xy} = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \sum_{i=1}^n x_i y_i - \frac{1}{n} \left(\sum_{i=1}^n x_i \right) \left(\sum_{i=1}^n y_i \right).$$

Naturligtvis har vi även t.ex. om s_x^2 är stickprovsvariansen för x -dataserien $S_{xx} = (n-1)s_x^2$.

2.4 Konfidensintervall för α och β

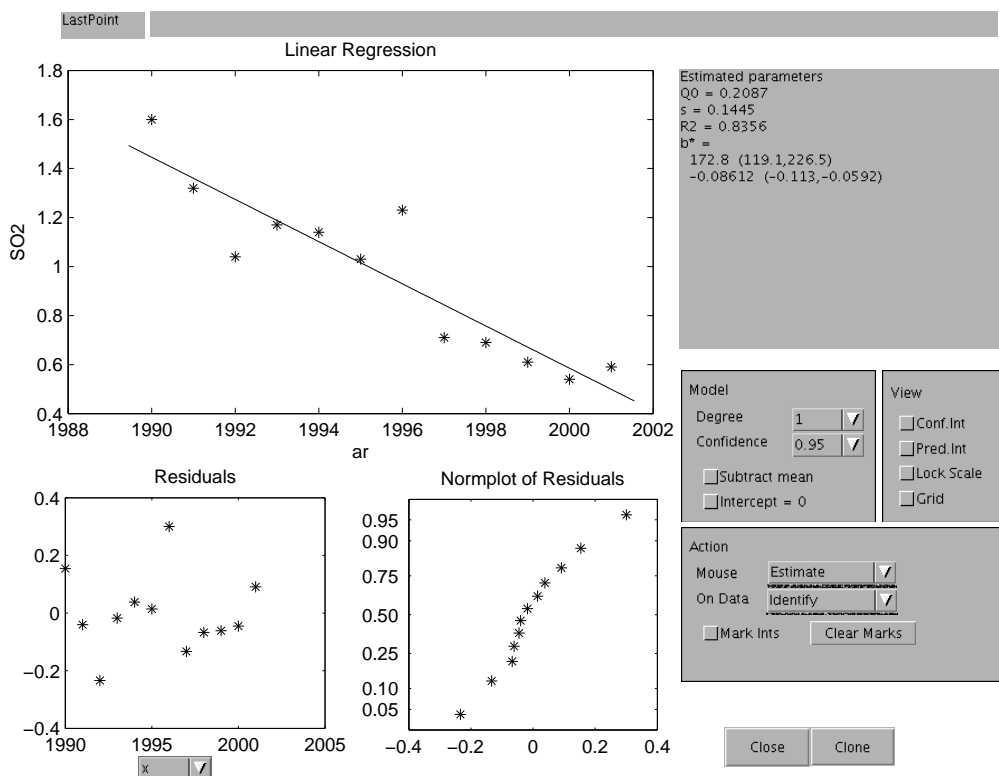
Eftersom skattningarna av α och β är normalfördelade får vi direkt konfidensintervall med konfidensgraden $1 - a$ (α är upptagen) precis som tidigare enligt

$$I_\beta = \beta^* \pm t_{a/2}(f) d(\beta^*) = \beta^* \pm t_{a/2}(n-2) \cdot \frac{s}{\sqrt{S_{xx}}}$$

$$I_\alpha = \alpha^* \pm t_{a/2}(f) d(\alpha^*) = \alpha^* \pm t_{a/2}(n-2) \cdot s \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}}$$

Om σ skulle råka vara känd används naturligtvis den i stället för s och då även λ - i stället för t -kvantiler.

Exempel 2.2. Hoburgsdata i exempel 2.1 analyserades, med hjälp av rutinen `reggui` i Matlab, och vi fick följande utskrift och figurer.



Figur 5: Regressionsanalys på materialet från Hoburgen; y = ”SO₂-halt” är plottad mot x = ”år”.

Överst till höger i utskriften ges en mängd information, bl.a. skattningar och konfidensintervall för modellens tre parametrar. För att göra det mer åskådligt sammanställer vi resultaten i en tabell:

parameter	skattning	95% konfidensintervall
α	172.8	(119.1, 226.5)
β	-0.08612	(-0.113, -0.0592)
σ	0.1445	

Vi ser att α skattas till 172.8 $\mu\text{g}/\text{m}^3$ och motsvarande intervall är $I_\alpha = (119.1, 226.5)$. Eftersom α är interceptet med y -axeln motsvaras α i detta exempel av SO_2 -halten vid år 0! Det går naturligtvis ej att anta att det linjära sambandet sträcker sig så långt bak, skattningen av α ger oss alltså inte omedelbart någon användbar information. Desto intressantare är lutningen β eftersom den talar om för oss hur mycket SO_2 -halten ändras under ett år. Från utskriften ser vi att denna förändring skattas till $-0.08612 \mu\text{g}/\text{m}^3$ per år. Intervallet $I_\beta = (-0.113, -0.0592)$ kan användas för att testa hypotesen $H_0 : \beta = 0$, vilket skulle innebära att SO_2 -halt inte påverkas av årtalen (d.v.s. ingen trend i data). Eftersom detta intervall inte täcker över 0 kan vi förkasta hypotesen $H_0 : \beta = 0$ och vi har påvisat (95% säkerhet) en nedåtgående *trend* i SO_2 -halt vid Hoburgen.

Vi ser också att σ skattas till 0.1445 (något konfidensintervall för denna storhet ges ej i utskriften). Residualkvadratsumman Q_0 är 0.2087 och det gäller som tidigare att $(\sigma^2)^* = s^2 = \frac{Q_0}{n-2}$ där n är antalet observerade talpar, d.v.s. $n=12$.

Storheten R2 i utskriften kommenteras nedan i avsnittet om förklaringsgraden. □

2.5 Skattning av punkt på linjen

För ett givet värde x_0 är Y 's väntevärde $E(Y(x_0)) = \alpha + \beta x_0 = \mu_0$, dvs en punkt på den teoretiska regressionslinjen. μ_0 skattas med motsvarande punkt på den skattade regressionslinjen som $\mu_0^* = \alpha^* + \beta^* x_0$. Vi ser direkt att skattningen är väntevärdesriktig samt att den måste vara normalfördelad (linjär funktion av två normalfördelade skattningar). Ett enkelt sätt att bestämma skattningens varians får vi om vi återigen utnyttjar att β^* och \bar{Y} är oberoende av varandra (men inte av α^*)

$$\begin{aligned} V(\mu_0^*) &= V(\alpha^* + \beta^* x_0) = [V(\alpha^* + \beta^* \bar{x}) + V(\beta^* (x_0 - \bar{x}))] = [V(\bar{Y}) + V(\beta^* (x_0 - \bar{x}))] = \\ &= V(\bar{Y}) + (x_0 - \bar{x})^2 V(\beta^*) = \frac{\sigma^2}{n} + (x_0 - \bar{x})^2 \frac{\sigma^2}{S_{xx}} = \sigma^2 \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right) \implies \\ \mu_0^* &\in N \left(\mu_0, \sigma \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}}} \right). \end{aligned}$$

Vi får således ett konfidensintervall för μ_0 med konfidensgraden $1 - a$ som

$$I_{\mu_0} = \mu_0^* \pm t_{a/2}(f) d(\mu_0^*) = \alpha^* + \beta^* x_0 \pm t_{a/2}(n-2) s \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}}}.$$

Exempel 2.3. Från exempel 1.1 på sid 3: I ett slumpmässigt urval av bilar avsattes y =”bensinförbrukning i stadskörning” som funktion av x =”vikt” i en linjär regressionsmodell $Y_i = \alpha + \beta x_i + \varepsilon_i$, $\varepsilon_i \in N(0, \sigma)$. Parametrarna skattas enligt resultaten i avsnitt 2.3 till $\alpha^* = 0.46$, $\beta^* = 0.0076$ samt $\sigma^* = 1.009$.

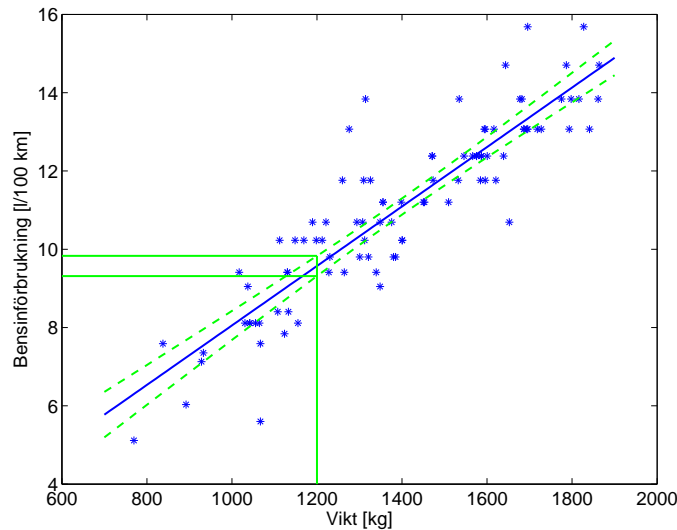
β är ett mått på hur mycket y beror av x , om vikten ökas med ett kg skattas ökningen av bensinförbrukningen med $\beta^* = 0.0076$ liter per 100 kilometer. Ett 95% konfidensintervall för β blir $I_\beta = (0.0068, 0.0084)$.

Antag att vi är speciellt intresserade av bilar som väger $x_0 = 1200$ kg. En skattning av medelförbrukningen μ_0 för denna typ av bilar blir då $\mu_0^* = \alpha^* + \beta^* x_0 = 9.57$ l/100 km. Ett

95% konfidensintervall för μ_0 blir med ovanstående uttryck $I_{\mu_0} = (9.32, 9.83]$. Detta intervall täcker alltså med sannolikhet 95% den sanna medelförbrukningen för bilar med vikt 1200 kg.

Observera att intervallet inte ger någon information om individuella 1200 kg bilars variation, så det är inte till så mycket hjälp till att ge någon uppfattning om en framtida observation (den 1200 kg bil du tänkte köpa?). Till detta behövs ett *prediktionsintervall*, se nästa avsnitt.

I figur 6 är konfidensintervallen förutom för 1200 kg bilar även plottat som funktion av vikten. I formeln för konfidensintervallet ser man att det är som smalast då $x_0 = \bar{x}$ vilket även kan antydast i figuren. Man ser även att observationerna i regel inte täcks av konfidensintervallen för linjen.



Figur 6: Bensinförbrukning enligt exempel 1.1. Skattad regressionslinje (—), konfidensintervall för linjen som funktion av vikt (- -). Konfidensintervall för linjen då vikten är $x_0 = 1200$ kg är markerat (-).

□

2.6 Prediktionsintervall för observationer

Intervallet ovan gäller väntevärdet för Y då $x = x_0$. Om man vill uttala sig om *en* framtida observation av Y för $x = x_0$ blir ovanstående intervall i regel för smalt. Om α , β och σ vore kända så skulle intervallet $\alpha + \beta x_0 \pm \lambda_{\alpha/2} \sigma$ täcka en framtida observation Y med sannolikhet $1 - \alpha$.

Eftersom regressionslinjen skattas med $\mu_0^* = \alpha^* + \beta^* x_0$ kan vi få hur mycket en framtida observation $Y(x_0)$ varierar kring den skattade linjen som

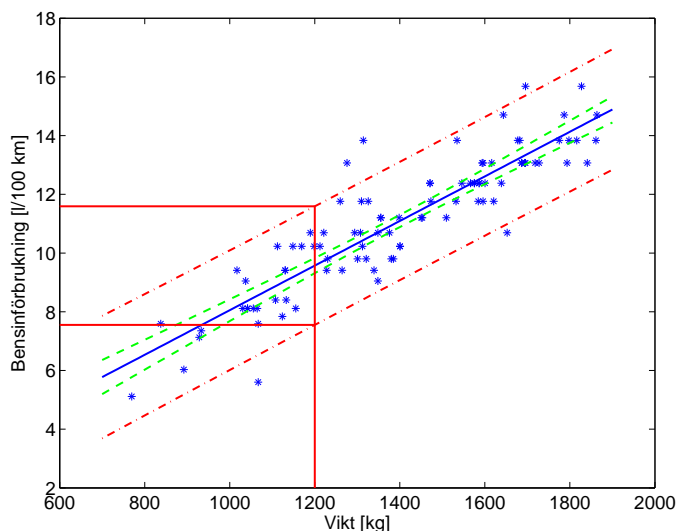
$$V(Y(x_0) - \alpha^* - \beta^* x_0) = V(Y(x_0)) + V(\alpha^* + \beta^* x_0) = \sigma^2 \left(1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right).$$

Vi kan alltså få ett *prediktionsintervall* med prediktionsgraden $1 - p$ för en framtida observation som

$$I_{Y(x_0)} = \alpha^* + \beta^* x_0 \pm t_{p/2}(n-2)s \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}}}.$$

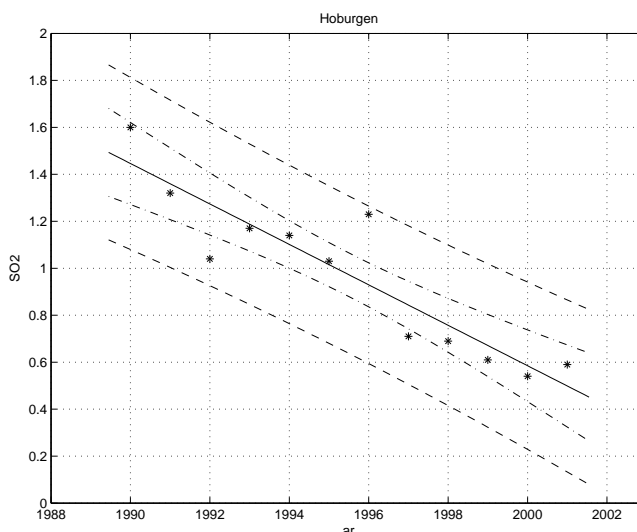
Observera att det bara är ettan i kvadratroten som skiljer mellan prediktionsintervallet och I_{μ_0} .

Exempel 2.4. Ett prediktionsintervall för bensinförbrukningen hos en 1200 kg bil enligt exempel 1.1 blir (7.6, 11.6) vilket är betydligt bredare än intervallet för väntevärdet. I figur 7 ses detta intervall och prediktionsintervallen som funktion av x_0 . □



Figur 7: Bensinförbrukning enligt exempel 1.1. Skattad regressionslinje (—), konfidensintervall för linjen som funktion av vikt (- -), prediktionsintervall för framtida observationer som funktion av vikt (- -). Prediktionsintervall för en framtida observation då vikten är $x_0 = 1200$ kg är markerat (-).

Exempel 2.5. Vi anknyter till exemplet med SO_2 -halterna igen. I figur 8 är både konfidensintervallet för linjens läge (det inre prick-streckade bandet) samt prediktionsintervallet (det yttre streckade bandet) uttridade som funktion av x_0 i Hoburgsdata.



Figur 8: Konfidensintervall för linjens läge (-) samt prediktionsintervall (- -) för SO_2 -halt ($\mu g/m^3$).

Vad är SO_2 -linjens läge vid år 1996, d.v.s vad är förväntad SO_2 -halt detta år? Ett 95% konfidensintervall för linjen beräknas till (0.83, 1.02) (jämför gärna med det inre bandet i figuren vid år 1996). Motsvarande prediktionsintervall (yttre band) för detta år är (0.59, 1.26), den uppmätta SO_2 -halten 1996 hade alltså, med 95% sannolikhet, kunnat hamna någonstans mellan 0.59 och 1.26 $\mu g/m^3$.

På motsvarande sätt kan man använda prediktionsintervallet för att säga att uppmätt SO_2 -halt år 2002, med 95% säkerhet, kommer att hamna någonstans i intervallet (0.03, 0.79) $\mu g/m^3$ (gör en försiktig extrapolation i figuren).

□

2.7 Kalibreringsintervall

Om man observerat ett värde y_0 på y , vad blir då x_0 ? Man kan lösa ut x_0 ur $y_0 = \alpha^* + \beta^* x_0$ och får

$$x_0^* = \frac{y_0 - \alpha^*}{\beta^*}$$

Denna skattning är inte normalfördelad, men vi kan t.ex använda Gauss approximationsformler för att få en skattning av $d(x_0^*)$ och konstruera ett approximativt intervall

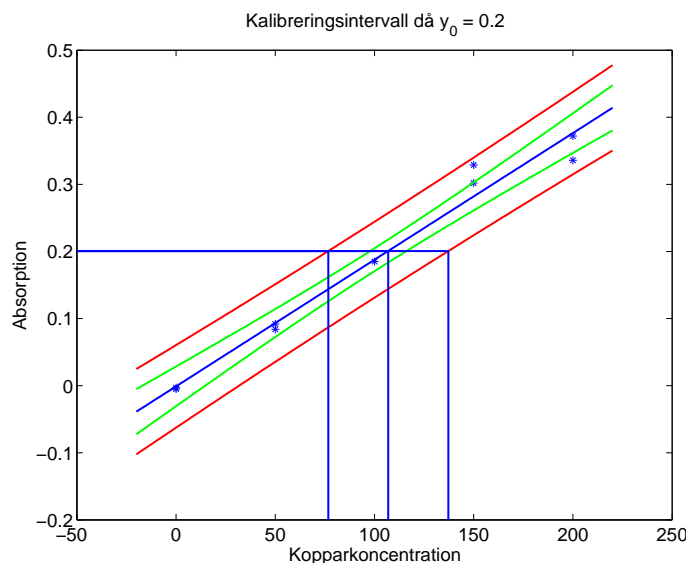
$$I_{x_0} = x_0^* \pm t_{a/2}(n-2)d(x_0^*) = \bar{x} + \frac{y_0 - \bar{y}}{\beta^*} \pm t_{a/2}(n-2) \cdot \frac{s}{|\beta^*|} \sqrt{1 + \frac{1}{n} + \frac{(y_0 - \bar{y})^2}{(\beta^*)^2 S_{xx}}}$$

Ett annat sätt att konstruera kalibreringsintervallet är att dra en linje $y = y_0$ och ta skärningspunkterna med prediktionsintervallet som gränser i kalibreringsintervallet. Ett analytiskt uttryck för detta blir efter lite arbete

$$I_{x_0} = \bar{x} + \frac{\beta^*(y_0 - \bar{y})}{c} \pm \frac{t_{p/2}(n-2) \cdot s}{c} \sqrt{c(1 + \frac{1}{n}) + \frac{(y_0 - \bar{y})^2}{S_{xx}}}$$

$$c = (\beta^*)^2 - \frac{(t_{p/2}(n-2) \cdot s)^2}{S_{xx}}$$

Uttrycket gäller då β är signifikant skild från noll annars är det inte säkert att linjen skär prediktionsintervallet. Grafiskt konstrueras detta intervall enligt figur 9.



Figur 9: Kalibreringsintervall konstruerat som skärning med prediktionsintervall. I försöket har man för ett par prover med kända kopparkoncentrationer mätt absorption med atomabsorptionsspektrofotometri. Kalibreringsintervallet täcker med ungefär 95% sannolikhet den rätta kopparkoncentrationen för ett prov med okänd kopparkoncentration där absorptionen uppmätts till 0.2.

2.8 Modellvalidering

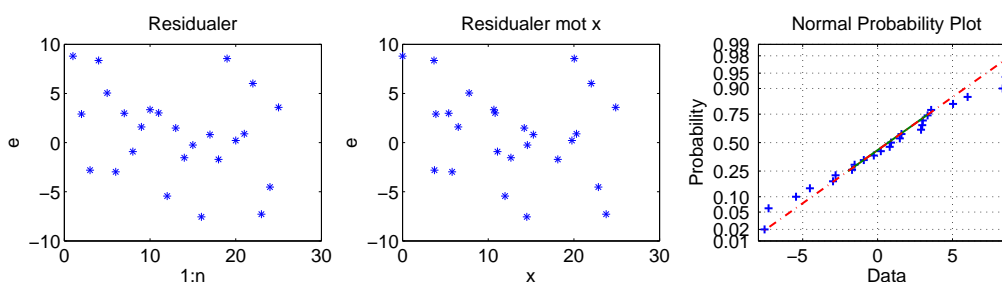
2.8.1 Residualanalys

Modellen vi använder baseras på att avvikelserna från regressionslinjen är likafördelade ($\varepsilon_i \in N(0, \sigma)$) och oberoende av varandra vilket medför att även observationerna Y_i är normalfördelade och oberoende. Dessa antaganden används då vi tar fram fördelningen för skattningarna. För att övertyga sig om att antagandena

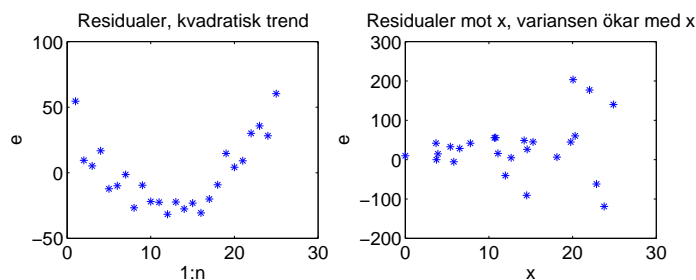
är rimliga kan det vara bra att studera avvikelserna mellan observerade y -värden och motsvarande punkt på den skattade linjen, d.v.s. de sedan tidigare definierade *residualerna*

$$r_i = y_i - (\alpha^* + \beta^* x_i), \quad i = 1, \dots, n,$$

eftersom dessa är observationer av ε_i . Residualerna bör alltså se ut att komma från en och samma normalfördelning samt vara oberoende av dels varandra, samt även av alla x_i . I figur 10 visas några exempel på residualplottar som ser bra ut medan de i figur 11 ser mindre bra ut.



Figur 10: Bra residualplottar. Residualerna plottade i den ordning de kommer, mot x samt i en normalfördelningsplott. De verkar kunna vara oberoende normalfördelade observationer.

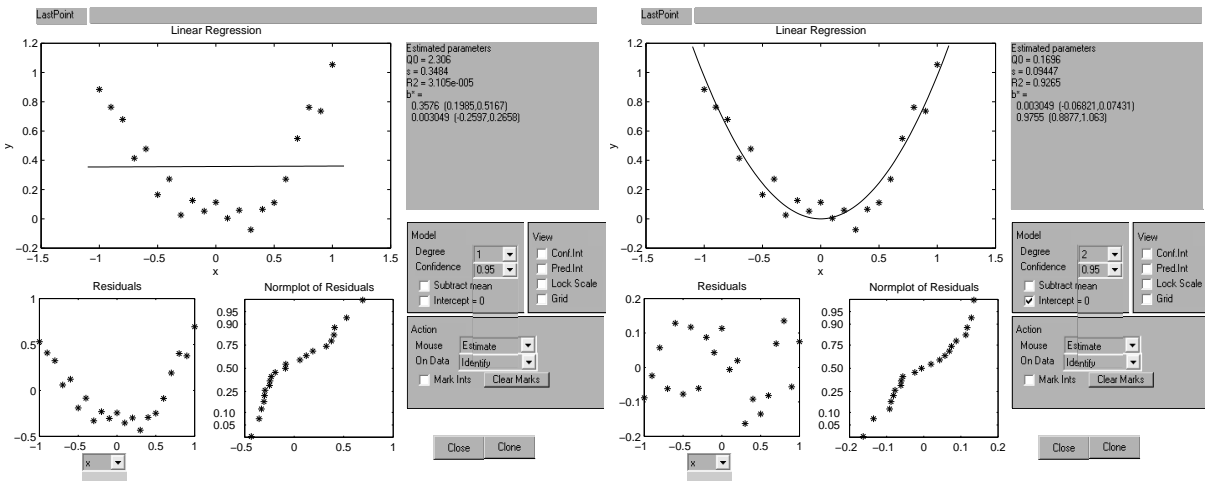


Figur 11: Residualplottar där man ser en tydlig kvadratisk trend i den vänstra figuren och i den högra ser man att variansen ökar med ökat x .

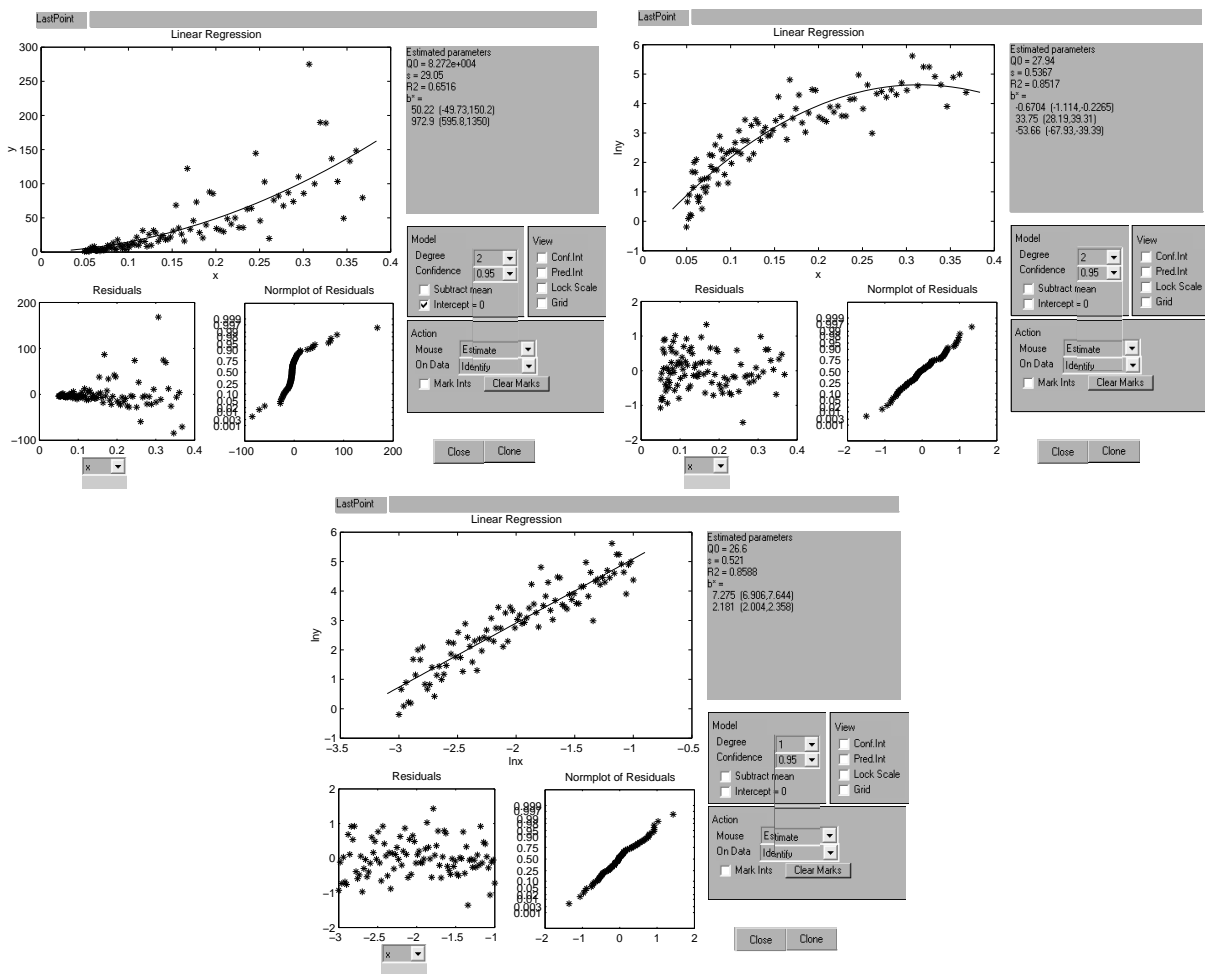
Exempel 2.6. Genom att studera graferna i figur 5 kan vi undersöka om den linjära modellen passar bra till Hoburgsdata. Residualplotten (nederst till vänster) visar inte några oroväckande trender och normalfördelningsplotten (nederst till höger) gör det rimligt att avvikelserna (residualerna) är normalfördelade. Sammantaget verkar det linjära modellen med oberoende och normalfördelningsantagande vara rimlig i detta fall. \square

Exempel 2.7. I figur 12a) anpassades modellen $y_i = \alpha + \beta x_i + \varepsilon_i$. Residualplotten i nedre vänstra hörnet säger att residualvärdet beror på x . Sambandet är alltså inte linjärt, snarare kvadratisk. Om vi istället anpassar modellen $y_i = \alpha + \beta_1 x_i + \beta_2 x_i^2 + \varepsilon_i$ ser residualerna ut som de ska (se figur 12b). \square

Exempel 2.8. Anpassa den kvadratiske modellen $y_i = \beta_1 x_i + \beta_2 x_i^2 + \varepsilon_i$ (se figur 13a). Anpassningen är dålig eftersom residualernas varians ökar med x . För att åtgärda det anpassar vi istället modellen $\ln y_i = \alpha + \beta_1 x_i + \beta_2 x_i^2 + \varepsilon_i$ (se figur 13b). Anpassningen är bättre eftersom residualvariansen nu är konstant. Däremot kan vi vara lite tveksamma till en kvadratisk modell eftersom modellen då säger att y ska avta för stora x . Det stämmer inte med observationerna. En bättre transformation är då att istället anpassa modellen $\ln y_i = \alpha + \beta_1 \ln x_i + \varepsilon_i$ (se figur 13c). Nu ser residualerna ut som de ska. \square

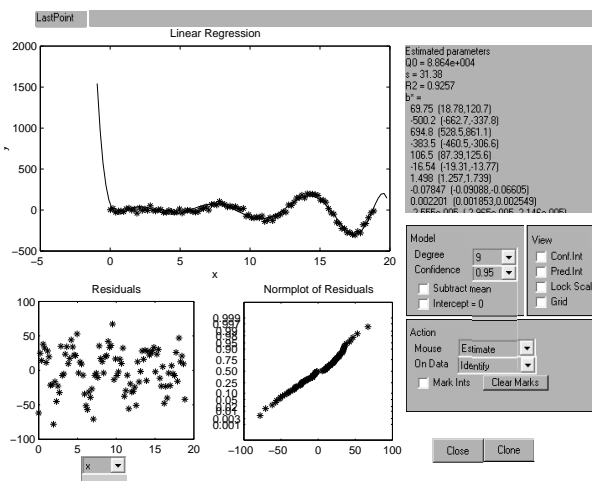


Figur 12: (a) Anpassning av linjär modell till kvadratiska data (vänster). (b) Anpassning av kvadratisk modell till kvadratiska data (höger).



Figur 13: (a) Anpassning av kvadratisk modell (överst till vänster) (b) Anpassning av kvadratisk modell efter logartimering av y (överst till höger) (c) Anpassning av linjär modell efter logartimering av både y och x (underst)

Exempel 2.9. Det är inte säkert att det går att hitta en linjär modell eller en enkel transformation som passar. Anpassa modellen $y_i = \alpha + \beta_1 x_i + \dots + \beta_p x_i^p + \varepsilon_i$ (se figur 14). Trots att vi anpassat ett polynom av högt gradtal finns det fortfarande struktur i residualerna och någon enkel transformation som skulle hjälpa är svårt att tänka ut! Antingen är det inte linjärt eller så är det inte oberoende, eller båda, kanske är det en tidsserie². Vill man lösa det problemet får man läsa *Stationära stokastiska processer*. \square



Figur 14: Anpassning av polynom till icke-linjärt samband

2.8.2 Är β signifikant?

Eftersom β anger hur mycket y beror av x är det även lämpligt att ha med följande hypotestest i en modellvalidering

$$H_0 : \beta = 0$$

$$H_1 : \beta \neq 0$$

t.ex. genom att förkasta H_0 om punkten 0 ej täcks av I_β . Om H_0 inte kan förkastas har y inget signifikant beroende av x och man kan kanske använda modellen $Y_i = \alpha + \varepsilon_i$ i stället.

2.9 Förklaringsgrad

En vanlig teknik när man analyserar data är att man försöker dela upp den variation som ses i mätningarna på olika variationskällor. Vid enkel linjär regression gäller uppdelningen:

”Total variation” = ”variation förklarad av linjen” + ”oförklarad variation”, där

- ”total variation” = $\sum_{i=1}^n (y_i - \bar{y})^2$, d.v.s. den variation som finns i y -värdena utan att vi tar hänsyn till x -värdena
- ”variation förklarad av linjen” = $\sum_{i=1}^n ((\alpha^* + \beta^* x_i) - \bar{y})^2$, vilket tolkas som den del av variationen i y -led som beskrivs av den linjära modellen
- ”oförklarad variation” = $\sum_{i=1}^n (y_i - (\alpha^* + \beta^* x_i))^2$, vilket är identiskt med residualkvadratsumman Q_0 och tolkas som den ”återstående” variation vi inte kan förklara med den linjära modellen.

²Modellen är i själva verket icke-linjär: $y_i = \sin(x_i) \cdot x_i^2 + \varepsilon_i$

Ett mått på hur väl linjen förklarar data är kvoten mellan variation förklarad av linjen och total variation. Denna kvot är *förklaringsgraden*

$$R^2 = \frac{\sum_{i=1}^n ((\alpha^* + \beta^* x_i) - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

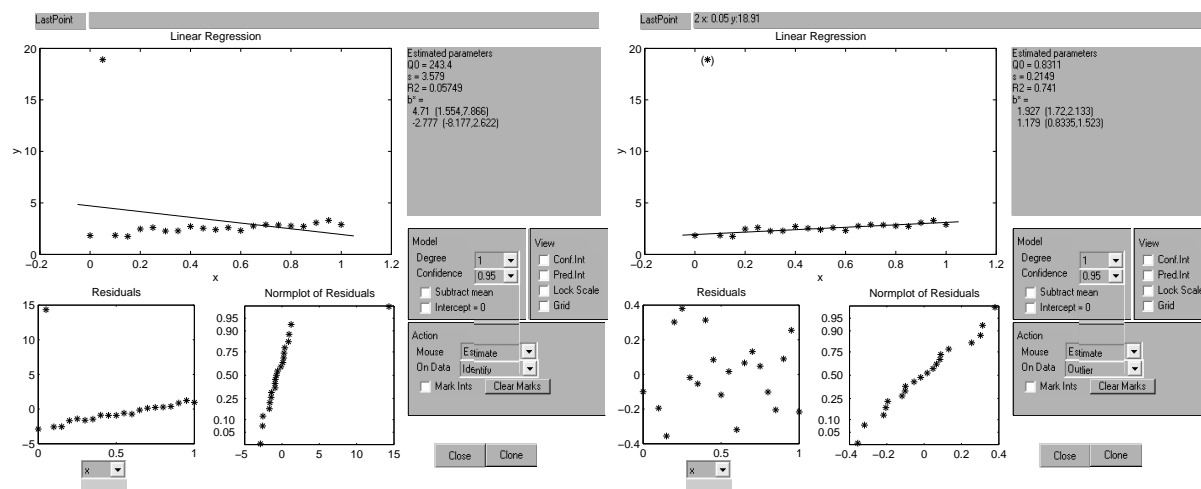
som ligger mellan noll och ett. Om R^2 har ett värde nära ett ligger talparen nära en rät linje - data kan därmed förklaras väl av den linjära modellen. Ett R^2 -värde nära noll tyder på att data ej uppvisar ett speciellt linjärt samband och därmed inte förklaras bra av vår linjära modell.

Exempel 2.10. Vid regressionsanalysen på Hoburgsdata i exempel 2.2. blev $R^2=0.8356$. Huvudparten, 84%, av den variation vi ser i SO_2 -halt kan alltså förklaras med den linjärt avtagande trenden i mätningarna. \square

Förklaringsgraden är identisk med kvadraten på korrelationskoefficienten, se avsnitt 4.

2.10 Outliers

Det är viktigt att vara uppmärksam på *outliers*, dvs enskilda observationer som ligger misstänkt långt från de övriga och som får ett stort inflytande på skattningen av linjen (se figur 15). Outliers kan vara rena felmätningar, i så fall bör de korrigeras eller plockas bort, men de kan också bero på naturlig variation i data. Då bör man överväga en modell som kan ta hänsyn till den variationen eller använda en mer robust skattningsmetod (ingår ej i denna kurs).



Figur 15: (a) Anpassad modell med en outlier (vänster) (b) Anpassad modell med outliern bortplockad (höger).

2.11 Linjärisering av några icke linjära samband

Vissa typer av exponential- och potenssamband med multiplikativa fel kan logaritmeras för att få en linjär relation. T.ex. fås när man logaritmerar

$$z_i = a \cdot e^{\beta x_i} \cdot \varepsilon'_i \quad \xrightarrow{\ln} \quad \underbrace{\ln z_i}_{y_i} = \underbrace{\ln a}_{\alpha} + \beta \cdot x_i + \underbrace{\ln \varepsilon'_i}_{\varepsilon_i}$$

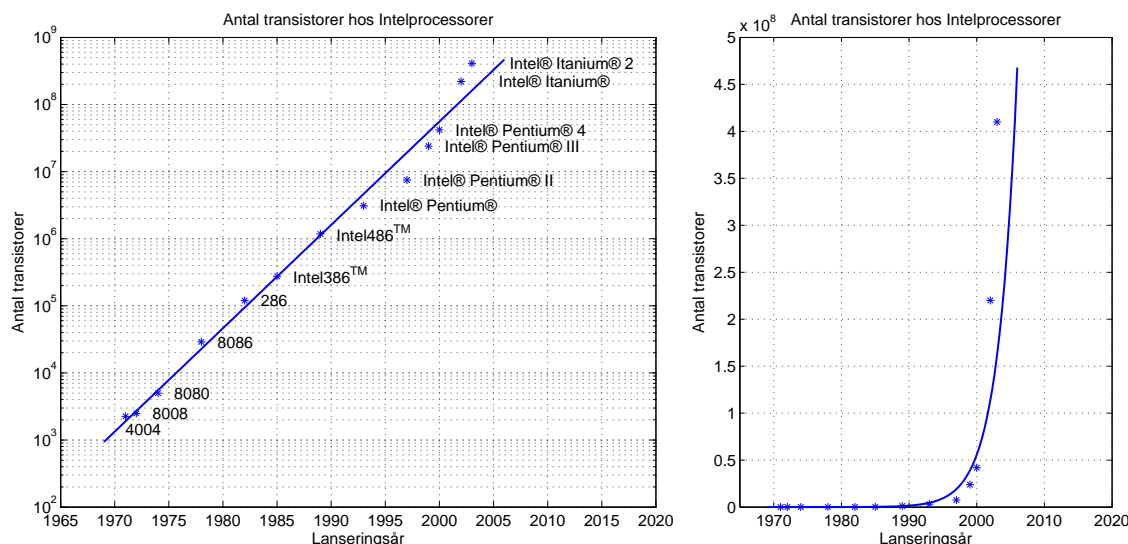
ett samband på formen $y_i = \alpha + \beta x_i + \varepsilon_i$. Man logaritmerar således z_i -värdena och skattar α och β som vanligt och transformerar till den ursprungliga modellen med $a^* = e^{\alpha^*}$. Observera att de multiplikativa felen

ε'_i bör vara lognormalfördelade (dvs $\ln \varepsilon'_i \in N(0, \sigma)$). En annan typ av samband är

$$z_i = a \cdot t_i^\beta \cdot \varepsilon'_i \quad \xrightarrow{\ln} \quad \underbrace{\ln z_i}_{y_i} = \underbrace{\ln a}_{\alpha} + \beta \underbrace{\ln t_i}_{x_i} + \underbrace{\ln \varepsilon'_i}_{\varepsilon_i}$$

där man får logaritmera både z_i och t_i för att få ett linjärt samband.

I figur 16 ses ett exempel där logaritmering av y -värdena ger ett linjärt samband.



Figur 16: Antal transistorer på en cpu mot lanseringsår med logaritmisk y -axel i vänstra figuren. Till höger visas samma sak i linjär skala. Det skattade sambandet är $y = 5.13 \cdot 10^{-301} \cdot e^{0.35x}$.

2.12 Jämförelse av två lutningar

Ibland har man en situation där man vill undersöka om regressionssambandet kan vara identiskt för olika grupper. Är t.ex. sambandet mellan blodtryck och ålder det samma för både män och kvinnor? Speciellt intressant kan det vara att studera om den årliga blodtrycksökningen är likartad för de båda könen. Om vi som modell använder två linjära regressionssamband (en för kvinnor och en för män) motsvaras problemet av att jämföra lutningarna i de två sambanden, d.v.s. undersöka om $\beta_{\text{kvinnor}} = \beta_{\text{män}}$. Ett exempel får illustrera metodiken.

Exempel 2.11. SO_2 -halten bestämdes inte enbart vid Hoburgen på Gotland utan även vid Rörvik i norra Halland (figur 17). Är trenden i SO_2 -halt den samma vid de två mätstationerna eller skiljer den sig åt?

Vi tänker oss att för Hoburgen och mätningarna $(x_1, y_1), \dots, (x_i, y_i), \dots, (x_{n_H}, y_{n_H})$ har vi modellen

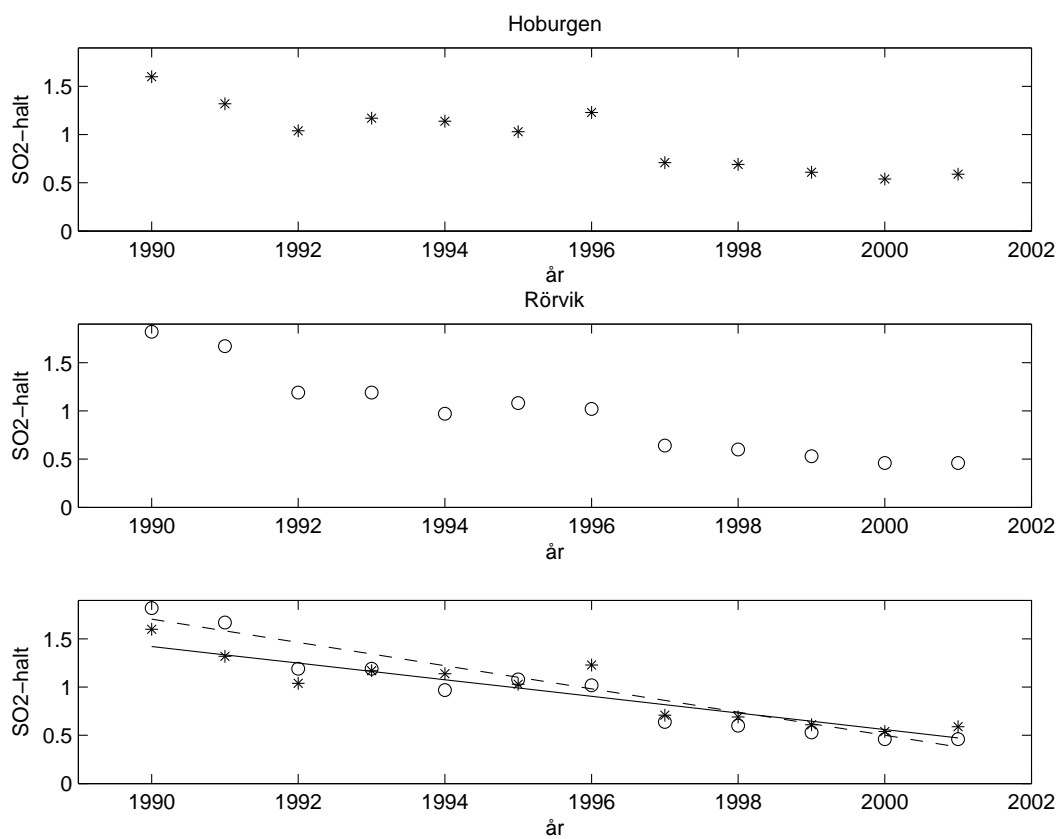
$$y_i = \alpha_H + \beta_H \cdot x_i + \varepsilon_i, \quad \varepsilon_i \in N(0, \sigma_H)$$

och för Rörvik och mätningarna $(x_1, y_1), \dots, (x_j, y_j), \dots, (x_{n_R}, y_{n_R})$ har vi modellen

$$y_j = \alpha_R + \beta_R \cdot x_j + \varepsilon_j, \quad \varepsilon_j \in N(0, \sigma_R).$$

Genom att göra två separata analyser i Matlab får vi för Hoburgen skattningarna (resultaten är hämtade från exempel 2.2).

$$\alpha_H^* = 172.8; \quad \beta_H^* = -0.08612; \quad \sigma_H^* = 0.1445$$



Figur 17: SO₂-halt vid Hoburgen (överst) samt i Rörvik (mitten). Underst visas mätningarna från båda stationerna med skattade regressionslinjer utritade (heldragen linje för Hoburgen och streckad för Rörvik)

medan motsvarande för Rörvik är

$$\alpha_R^* = 241.5; \beta_R^* = -0.1205; \sigma_R^* = 0.1436$$

Nu är vi intresserade av hur stor $\beta_R - \beta_H$ är och en **skattning** av denna storhet kan vi få genom $\beta_R^* - \beta_H^* = -0.1205 - (-0.08612) = -0.0344$.

Vill vi göra **konfidensintervall** för differensen $\beta_R - \beta_H$ måste vi ha en uppfattning om "hur bra" denna skattning är, d.v.s. veta variansen för $\beta_R^* - \beta_H^*$. Men från tidigare vet vi att

$$V(\beta_R^*) = \frac{\sigma_R^2}{S_{Rxx}}$$

där $S_{Rxx} = \sum(x_j - \bar{x})^2$ är kvadratsumman på de x -värden som användes vid Rörviksmätningarna. För Hoburgen har vi på motsvarande sätt

$$V(\beta_H^*) = \frac{\sigma_H^2}{S_{Hxx}}$$

där S_{Hxx} är kvadratsumman på de x -värden som användes vid Hoburgsmätningarna. Men eftersom x -värdena består av 11 årtal med start 1990 och slut 2001 och vi dessutom mäter vid samma år vid de två stationerna gäller att $S_{Hxx} = S_{Rxx} = 143$.

Om vi dessutom kan anta att $\sigma_H = \sigma_R$ (verkar rimligt i detta exempel) kan vi kalla denna gemensamma standardavvikelse för σ . Detta ger

$$V(\beta_R^* - \beta_H^*) = V(\beta_R^*) + V(\beta_H^*) = \sigma^2 \left(\frac{1}{S_{Rxx}} + \frac{1}{S_{Hxx}} \right).$$

För att beräkna en skattning av den gemensamma standardavvikelsen gör vi en "poolning" av standardavvikelserna av samma slag som tidigare (observera $n-2$)

$$\begin{aligned} \sigma^{2*} &= \frac{(n_R - 2) \cdot \sigma_R^{2*} + (n_H - 2) \cdot \sigma_H^{2*}}{(n_R - 2) + (n_H - 2)} = \\ &= \frac{(12 - 2) \cdot 0.1436^2 + (12 - 2) \cdot 0.1445^2}{(12 - 2) + (12 - 2)} = 0.0208. \end{aligned}$$

Nu kan vi konstruera ett 95% intervall på välbekant sätt:

$$\begin{aligned} I_{\beta_R - \beta_H} &= (\beta_R^* - \beta_H^* \pm t_{\alpha/2}(n_R - 2 + n_H - 2)d(\beta_R^* - \beta_H^*)) = \\ &= (\beta_R^* - \beta_H^* \pm t_{\alpha/2}(n_R - 2 + n_H - 2)\sqrt{\sigma^{2*}\left(\frac{1}{S_{Rxx}} + \frac{1}{S_{Hxx}}\right)}) = \\ &= (-0.0344 \pm 2.09 \cdot \sqrt{0.0208\left(\frac{1}{143} + \frac{1}{143}\right)}) = (-0.0344 \pm 0.0356) = (-0.070, 0.0012). \end{aligned}$$

Eftersom detta intervall täcker över 0 har vi inte påvisat att det finns en skillnad mellan lutningarna. Dessa mätningar tyder alltså inte på att trenden i SO_2 skiljer sig åt vid de två stationerna. \square

3 Multipel linjär regression på matrisform

Med matrisnotation kan en allmän linjär regressionsmodell med p st förklarande x -variabler, av typen

$$y_i = \alpha + \beta_1 x_{1i} + \dots + \beta_p x_{pi} + \varepsilon_i$$

vare sig den är enkel eller multipel, skrivs

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e},$$

där de ingående matriserna har följande form:

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}, \quad \mathbf{X} = \begin{pmatrix} 1 & x_{11} & \dots & x_{p1} \\ 1 & x_{12} & \dots & x_{p2} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{1n} & \dots & x_{pn} \end{pmatrix}, \quad \boldsymbol{\beta} = \begin{pmatrix} \alpha \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix} \quad \text{och} \quad \mathbf{e} = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}.$$

Rent allmänt fås minsta-kvadratlösningen $\boldsymbol{\beta}^*$ till ett överbestämt ekvationssystem $\mathbf{y} = \mathbf{X}\boldsymbol{\beta}$ via de så kallade normalekvationerna

$$\mathbf{X}^t \mathbf{X} \boldsymbol{\beta} = \mathbf{X}^t \mathbf{y},$$

som $\boldsymbol{\beta}^* = (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \mathbf{y}$. Man bör dock i möjligaste mån undvika att lösa ut $\boldsymbol{\beta}$ genom att invertera matrisen $\mathbf{X}^t \mathbf{X}$. Om matrisen är illa konditionerad kan man nämligen få en feltillväxt som gör resultatet helt oanvändbart. En numeriskt sett effektivare och mer stabil lösning fås om man i Matlab använder operatoren `\` som kan uppfattas som vänsterdivision.

Det rekommenderade sättet att lösa matrisekvationen ovan är alltså

```
>> b = X \ y
```

Skattningen av σ fås genom

$$\sigma^* = s = \sqrt{\frac{Q_0}{n - (p + 1)}}$$

där Q_0 kan beräknas antingen som $Q_0 = \mathbf{y}^t \mathbf{y} - \boldsymbol{\beta}^{*t} \mathbf{X}^t \mathbf{y}$, eller genom att utnyttja att $Q_0 = \sum_{i=1}^n r_i^2 = \mathbf{r}^t \mathbf{r}$ där

residualerna $\mathbf{r} = \mathbf{y} - \mathbf{X}\boldsymbol{\beta}$. Den s.k. *kovariansmatrisen* för $\boldsymbol{\beta}^*$ ges av $\sigma^2 (\mathbf{X}^t \mathbf{X})^{-1}$ vilket innebär att medelfelen $\mathbf{d}(\beta_0^*)$, $\mathbf{d}(\beta_1^*)$, etc, fås som roten ur respektive diagonalelement i $s^2 (\mathbf{X}^t \mathbf{X})^{-1}$. Den skattade linjen i punkten

$\mathbf{x}_0 = (1 \ x_0^{(1)} \ x_0^{(2)})$ ges av $\mu_0^* = \mathbf{x}_0 \boldsymbol{\beta}^* \in N\left(\mu_0, \sigma \sqrt{\mathbf{x}_0 (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{x}_0^t}\right)$.

Exempel 3.1. För att undersöka pressningstemperaturens och pressningstryckets inverkan vid tillverkning av en typ av plastkomposit iordningställdes två provbitar för var och en av fem kombinationer av tryck och temperatur. Böjspänningen hos de olika provbitarna av plastkompositen mättes och blev

Böjspänning (y) (N/mm ²)	Temperatur (x_1) (°C)	Tryck (x_2) (kg/cm ²)
152	180	450
150	180	450
103	190	375
99	190	375
88	200	350
89	200	350
122	210	375
120	210	375
162	220	450
161	220	450

Anpassa modellen $y_i = \alpha + \beta_1 x_{1i} + \beta_2 x_{2i} + \varepsilon_i$ och gör ett 95 % konfidensintervall för hur mycket böjspänningen ökar då temperaturen ökar med 1 °C. Gör också ett 95 % prediktionsintervall för böjspänningen då temperaturen är 200 °C och trycket 400 kg/cm².

Lösning: Skriv om modellen $y_i = \alpha + \beta_1 x_{1i} + \beta_2 x_{2i} + \varepsilon_i$ som $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$ med

$$\mathbf{y} = \begin{pmatrix} 152 \\ 150 \\ 103 \\ 99 \\ 88 \\ 89 \\ 122 \\ 120 \\ 162 \\ 161 \end{pmatrix} \quad \mathbf{X} = \begin{pmatrix} 1 & 180 & 450 \\ 1 & 180 & 450 \\ 1 & 190 & 375 \\ 1 & 190 & 375 \\ 1 & 200 & 350 \\ 1 & 200 & 350 \\ 1 & 210 & 375 \\ 1 & 210 & 375 \\ 1 & 220 & 450 \\ 1 & 220 & 450 \end{pmatrix}, \quad \boldsymbol{\beta} = \begin{pmatrix} \alpha \\ \beta_1 \\ \beta_2 \end{pmatrix}$$

Parameterskattningar blir

$$\boldsymbol{\beta}^* = \mathbf{X} \backslash \mathbf{y} = \begin{pmatrix} -215.7 \\ 0.41 \\ 0.65 \end{pmatrix} = \begin{pmatrix} \alpha^* \\ \beta_1^* \\ \beta_2^* \end{pmatrix}$$

och, eftersom $Q_0 = \mathbf{r}^t \mathbf{r} = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}^*)^t (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}^*) = 243.63$,

$$\sigma^* = s = \sqrt{\frac{Q_0}{n - (p + 1)}} = \sqrt{\frac{243.63}{10 - (2 + 1)}} = 5.90.$$

Ökningen i böjspänning då temperaturen ökar en grad ges av β_1 . För att kunna beräkna konfidensintervall för β_1 behöver vi också beräkna

$$(\mathbf{X}^t \mathbf{X})^{-1} = \begin{pmatrix} 29.24 & -0.1 & -0.0229 \\ -0.1 & 0.0005 & 0 \\ -0.0229 & 0 & 0.0001 \end{pmatrix}$$

Sedan kan vi få medelfelet $\mathbf{d}(\beta_1^*) = s\sqrt{0.0005}$, där vi tagit andra diagonalelementet i $(\mathbf{X}^t \mathbf{X})^{-1}$. Det första diagonalelementet gäller ju α^* och det tredje β_2^* . Ett konfidensintervall för β_1 med konfidensgrad $1 - \alpha$ fås sedan på vanligt sätt som

$$\begin{aligned} I_{\beta_1} &= (\beta_1^* \pm t_{\alpha/2}(n - (p + 1)) \cdot \mathbf{d}(\beta_1^*)) \\ &= (0.41 \pm \underbrace{t_{0.025}(7)}_{2.36} \cdot 5.90\sqrt{0.0005}) \\ &= (0.098, 0.722) \text{ N/mm}^2 \text{ per } ^\circ\text{C}. \end{aligned}$$

För att göra ett prediktionsintervall för Y_0 då $x_0^{(1)} = 200$ °C och $x_0^{(2)} = 400$ kg/cm² sätter vi $\mathbf{x}_0 = (1 \ 200 \ 400)$ och får skattningen av sambandet till $\mu_0^* = \mathbf{x}_0 \boldsymbol{\beta}^* = 124.6$ med medelfelet $\mathbf{d}(\mu_0^*) = s\sqrt{\mathbf{x}_0 (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{x}_0^t} = 5.90\sqrt{0.1} = 6.187$. Eftersom vi vill ha ett prediktionsintervall, inte ett konfidensintervall, ska vi lägga till en etta under rottecknet så att intervallet ges av

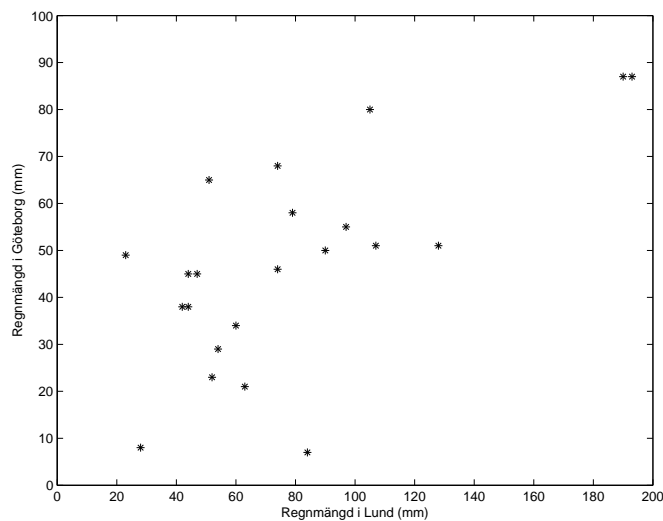
$$\begin{aligned} I_{Y(\mathbf{x}_0)} &= \left(\mathbf{x}_0 \boldsymbol{\beta}^* \pm t_{\alpha/2}(n - (p + 1)) \cdot s\sqrt{1 + \mathbf{x}_0^t (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{x}_0} \right) \\ &= (124.6 \pm \underbrace{t_{0.025}(7)}_{2.36} \cdot 5.90\sqrt{1 + 0.1}) = (110.0, 139.2) \text{ N/mm}^2 \end{aligned}$$

□

4 Korrelationsanalys

Regressionsanalysen i föregående avsnitt förutsatte att x -variablerna var ”fixa” i den meningen att de var uppmätta med inget eller försumbart mätfel. Om detta inte är uppfyllt är det lämpligare att göra en korrelationsanalys där man inte försöker anpassa någon regressionsfunktion till data utan enbart mäter graden av samband.

Exempel 4.1. I exempel 1.2 på sidan 3 noterades månadsnederbörden, d.v.s. den totala mängden nederbörd (mm) under en månad, i Göteborg och Lund under åren 2005 och 2006. I figur 18 markerar varje punkt en månad där Göteborgs nederbörd avläses på y -axeln och Lunds på x -axeln.



Figur 18: Månadsvisa mätningar av nederbörden (mm) där y = ”nederbörd i Göteborg” är plottad mot x = ”nederbörd i Lund”.

□

Från figuren tycks det finnas ett positivt samband mellan nederbördsmätningarna från de två städerna - regnar det mycket en månad i den ena staden tenderar det också att göra det i den andra.

4.1 Mått på samband

Som ett mått på samband mellan två variabler X och Y används kovariansen eller korrelationskoefficienten mellan variablerna. Kovariansen definieras som

$$C(X, Y) = E[(X - \mu_x)(Y - \mu_y)],$$

där μ_x och μ_y är väntevärdena för X och Y . Korrelationskoefficienten, ρ_{xy} är den normerade storheten

$$\rho_{xy} = \frac{C(X, Y)}{D(X) \cdot D(Y)},$$

där $D(X) = \sqrt{V(X)}$ är standardavvikelsen för X (och motsvarande för $D(Y)$). För korrelationskoefficienten gäller alltid att $-1 \leq \rho_{xy} \leq 1$.

Tolkning av de två storheterna är oftast enklast då man betraktar motsvarande skattningar. Antag att vi har n mätningar vardera av de två variablerna och därmed de n talparen $(x_1, y_1), \dots, (x_n, y_n)$. En skattning av kovariansen är då

$$c_{xy} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}).$$

och av korrelationskoefficienten

$$\rho_{xy}^* = r_{xy} = \frac{c_{xy}}{s_x s_y} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}.$$

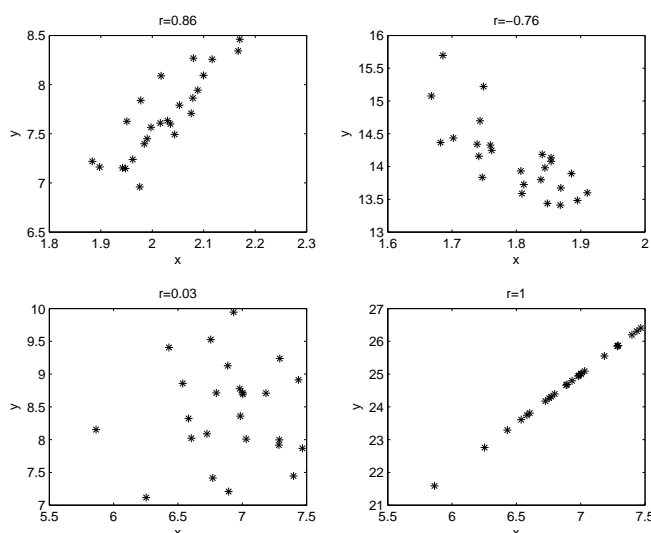
Observera att uttrycket $n - 1$ förkortats bort i sista ledet. För r_{xy} gäller att

$$-1 \leq r_{xy} \leq 1$$

samt att om vi har

positiv samvariation (positiv korrelation) mellan X och Y, d.v.s. $\rho_{xy} > 0$ tenderar $r_{xy} > 0$
 negativ samvariation (negativ korrelation) mellan X och Y, d.v.s. $\rho_{xy} < 0$ tenderar $r_{xy} < 0$
 ingen samvariation (ingen korrelation) mellan X och Y, d.v.s. $\rho_{xy} = 0$ tenderar $r_{xy} \approx 0$

Om $r_{xy} = 1$ innebär det att x -värdena och y -värdena ligger på en linje med positiv lutning; se figur 19.



Figur 19: Figureerna visar olika grad av samband med tillhörande korrelationskoefficient.

Observera att om r_{xy} ligger nära 0 tyder det på att det inte finns någon samvariation mellan de två variablerna (de är okorrelerade), däremot följer det inte att x och y är oberoende. Om x -värdena och y -värdena däremot är hämtade från normalfördelning är okorrelerad identiskt med oberoende.

4.2 Test av samband

I exemplet med månadsnederbörd från Lund och Göteborg gav beräkningar i Matlab att $r_{xy} = 0.662$. Data tyder alltså på en positiv samvariation - men är värdet på r_{xy} tillräckligt stort för att vi ska kunna tro på att det **verkligen** finns en samvariation och att det observerade resultatet inte bara är ett utslag av slumpen?

Om r_{xy} är en skattning av den korrelation, ρ_{xy} , som finns mellan de s.v. X och Y vill vi alltså undersöka om ρ_{xy} är 0. De intressanta hypoteserna är:

$$H_0 : \rho_{xy} = 0 \text{ (inget samband); } H_1 : \rho_{xy} \neq 0 \text{ (samband).}$$

För att testa detta används storheten

$$t = r_{xy} \sqrt{(n-2)/(1-r_{xy}^2)}.$$

Om data kommer från en bivariat normalfördelning gäller nämligen att t är t -fördelad med $n-2$ frihetsgrader när H_0 är sann.

Exempel 4.2. Med ett värde $r_{xy} = 0.662$ i nederbördsdata blir

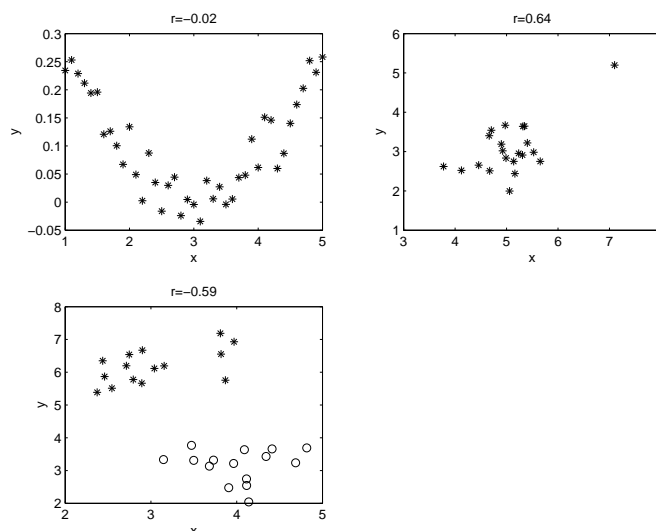
$$t = r_{xy} \sqrt{(n-2)/(1-r_{xy}^2)} = 0.662 \sqrt{(23-2)/(1-0.662^2)} = 3.95.$$

Eftersom 3.95 överstiger $t_{0.0005}(21) = 3.82$ innebär det att korrelationen är signifikant skild från 0 på nivå 0.001. Det finns alltså en positiv samvariation mellan de två städernas månadsnederbörd. \square

4.3 Var försiktig med korrelationskoefficienten!

Det finns en rad ”fallgropar” när man hanterar korrelationskoefficienter. Några exempel:

- r_{xy} mäter graden av linjärt samband - i figur 20(a) fås ett värde på r_{xy} som är ungefär 0 eftersom den negativa lutningen i figurens vänstra halva ”tas ut” av den positiva lutningen i andra halvan.
- r_{xy} är känslig för outliers, d.v.s. kraftigt avvikande värden kan starkt påverka värdet på korrelationskoefficienten. Utan outliern i figur 20(b) är $r_{xy}=0.24$, med outliern blir $r_{xy}=0.64$.
- r_{xy} kan bli missvisande då den används på mätningar som naturligt kan delas upp i två grupper (t.ex. kön) och där genomsnittsvärdena för x och y är olika i de två grupperna. I figur 20(c) verkar det inte finnas någon samvariation inom respektive grupp (eller eventuellt en positiv samvariation för ”stjärnorna”) men betraktar man hela materialet - och beräknar okritiskt r_{xy} - tyder korrelationskoefficienten på en negativ samvariation mellan X och Y .



Figur 20: Figurerna visar några situationer där korrelationskoefficienten inte okritiskt kan användas.

Samtliga dessa ”fällor” kan man förmodligen upptäcka om man alltid tar för vana att plotta sina data och inte bara slentrianmässigt beräknar korrelationskoefficienten.

Viktigare är det att komma ihåg att med korrelationskoefficienten mäter vi (och eventuellt påvisar) ett **statistiskt samband**. Det är därmed inte sagt att det finns ett **orsakssamband** mellan variablerna!

Exempel 4.3. Om man för ett antal städer noterar dels antal läkare i staden och dels antalet sjukdagar som stadens invånare har under ett år kommer man säkert att finna ett positivt samband mellan de två variablerna. Innebär det då att ju fler läkare man har i en stad medför det fler sjukdagar och att vi kan minska antalet sjukdagar genom att minska antalet läkare? Nej, naturligtvis inte; här är det en tredje faktor - antalet invånare i staden - som påverkar de båda undersökta variablerna. \square

4.4 Anknytning till linjär regression

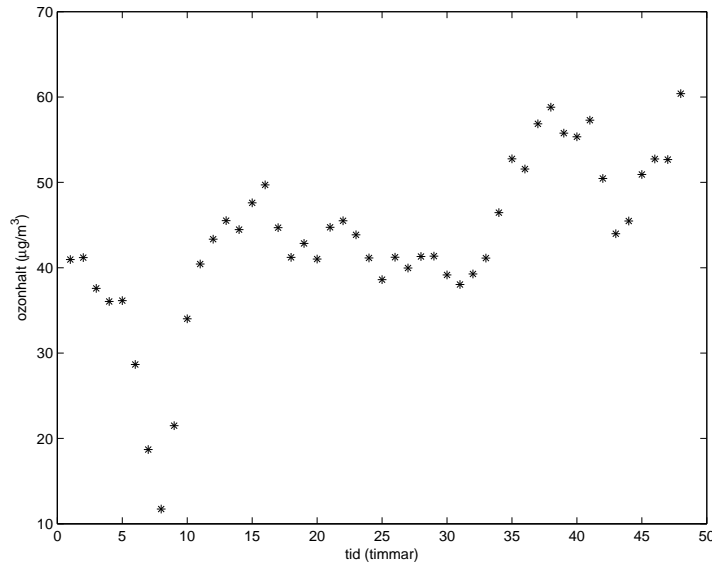
Korrelationskoefficienten mäter det *linjära* sambandet mellan x och y - alltså borde det kunna användas även vid linjär regression. I själva verket är kvadraten på korrelationskoefficienten matematiskt identisk med förklaringsgraden som beskrevs i avsnitt 2.9, d.v.s.

$$r_{xy}^2 = R^2.$$

Vid en regressionsanalys - antingen den beskrivs i datorprogram eller i rapporter - anges därför även ofta korrelationskoefficienten. Den är då ett mått på hur "stor nytta" man har av x -variabeln då man vill förutsäga y . Om r_{xy} är nära 1 (eller -1) betyder det att x och y ligger nästan på en linje och därmed kan y nästan förutsägas direkt utifrån x -värdet. Förklaringsgraden R^2 är då också nära 1. Om däremot värdet på r_{xy} är lågt (vilket ger en låg förklaringsgrad) är sambandet mellan variablerna svagt och y kan näppeligen förutsägas av enbart x .

Test av samband, som beskrivs i avsnitt 4.2, visar sig också vara identiskt med att testa att lutningen $\beta = 0$ (se avsnitt 2.8.2) i regressionsmodellen.

Observera dock - vilket vi redan påpekat - att det finns en skillnad i antagandena om x -värdena när det gäller regressionsanalys respektive korrelationsanalys. För förklaringsgraden R^2 i regressionsanalysen anses x -värdena vara fixa och att vi, i stort sett, kan själva bestämma dess värde. I korrelationsanalysen är däremot x -värdena och y -värdena "utbytbara".



Figur 21: Ozonhalt ($\mu\text{g}/\text{m}^3$) i Lund under 9 och 10 oktober 2001.

5 Tidsserier

En tidsserie är en uppsättning mätningar gjorda i tidsföljd, vi har mätningar y_1, y_2, \dots, y_n vid tidpunkterna $t = 1, 2, \dots, n$. Tidsserien skrivs ibland $\{y_t\}$.

Mätningarna är observationer av slumpvariablerna Y_1, Y_2, \dots, Y_n som utgör en process $\{Y_t\}$, ofta sägs tidsserien $\{y_t\}$ vara en realisering av processen $\{Y_t\}$. Det nya är att slumpvariablerna Y_1, Y_2, \dots, Y_n nu kan få vara beroende!

I avsnitt 2, se t.ex. exempel 2.1, studerade vi regressionsmodeller där x-variabeln var tiden, mätningarna $\{y_t\}$ var alltså en tidsserie. Då var vi noga med att avvikelserna från linjen, de stokastiska variablerna $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n$ var oberoende. Detta krav släpper vi alltså nu och tillåter dem att få vara beroende.

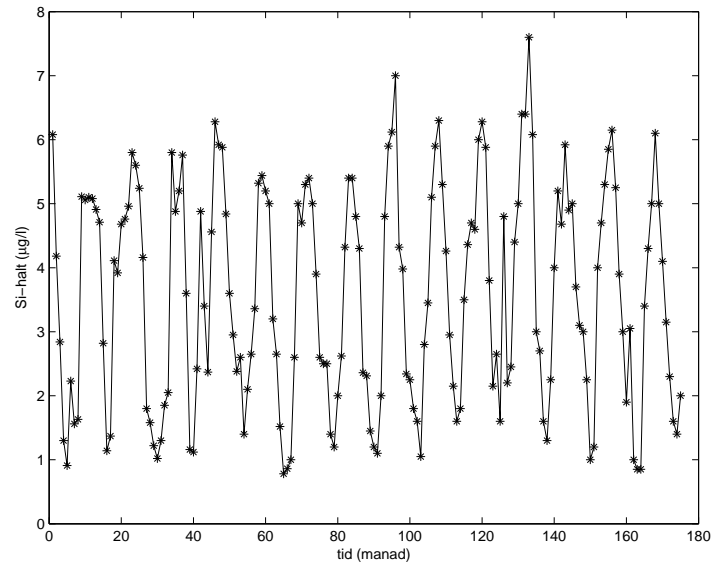
Exempel 5.1. I figur 21 visas halterna av ozon ($\mu\text{g}/\text{m}^3$) i Lund under 9 och 10 oktober 2001 (källa: www.oresundsluft.com/). Om man anpassar en linjär regressionsmodell till data kommer residualerna att visa ett tydligt mönster - de är kraftigt beroende (se kommande exempel). \square

Exempel 5.2. I figur 22 visas mätningar av kiselhalten ($\mu\text{g}/\text{l}$) i Ljungbyån under dryga 15 år (källa: SLU:s datavärdsbank, info1.ma.slu.se/db.html). För tydlighets skull är linjer dragna mellan punkterna. Mätningarna, som är gjorda en gång i månaden, visar ett tydligt säsongsmönster. \square

5.1 Syftet med analysen

I exemplen ovan ser vi några egenskaper som ofta karakteriserar tidsserier: en mjukt varierande trend, säsongsmönster och en "tröghet" som tyder på samband mellan i tiden näraliggande mätningar. Ett första steg i analysen är att **beskriva** tidsserien (grafiskt eller med numeriska mått) för att fånga upp dess karakteristiska drag. En grundläggande idé är att försöka dela upp $\{y_t\}$ i olika komponenter som t.ex. kan beskriva trend, säsong och kvarvarande "brus". En annan är att försöka beskriva beroendet mellan näraliggande mätningar, d.v.s. tidsseriens autokorrelation.

Mer avancerat är att försöka **modellera** komponenterna. Man talar t.ex. om *deterministiska modeller* där processens uppträdande är helt given enligt någon matematisk funktion t.ex. säsong = $A \cdot \sin(\varphi t)$. Mer fruktbart är oftast att använda en *slumpmodell* där det finns en slumpmässig variation mellan observationerna, t.ex. säsong = $A \cdot \sin(\varphi t) + e_t$, där e_t är en stokastisk variabel.



Figur 22: Månadsvisa mätningar av kiselhalt ($\mu\text{g}/\text{l}$) i Ljunbyån under dryga 15 år.

Ett vanligt syfte vid tidsserieanalys är att man vill kunna **förutsäga** kommande värden på den uppmätta variabeln. I miljösammanhang studerar man ofta tidsserier för att **övervaka** och upptäcka förändringar av en miljövariabel.

Analys av tidsserier är komplicerat, vill man fördjupa sig i det kan man läsa specialkurser i området. Vi kommer här endast att ge en orientering i ämnet där betoningen ligger på beskrivning av tidsserien samt på en enkel modell för autokorrelationen.

5.2 Beskrivning av tidsserien

Vi har tidigare studerat genomsnittsvärdet, väntevärdet μ , för en stokastisk variabel Y . Nu har vi en hel följd av stokastiska variabler $\{Y_t\}$, och motsvarande genomsnittsvärde, $E(Y_t) = \mu_t$, är en funktion av t och kan tolkas som trenden i tidsserien. Om genomsnittsfunktionen μ_t är konstant μ säges tidsserien vara *stationär*.

5.2.1 Komponentuppdelning

En bärande idé i tidsserieanalys är att försöka dela upp $\{y_t\}$ i flera olika komponenter. I sin enklaste form kan uppdelningen beskrivas som

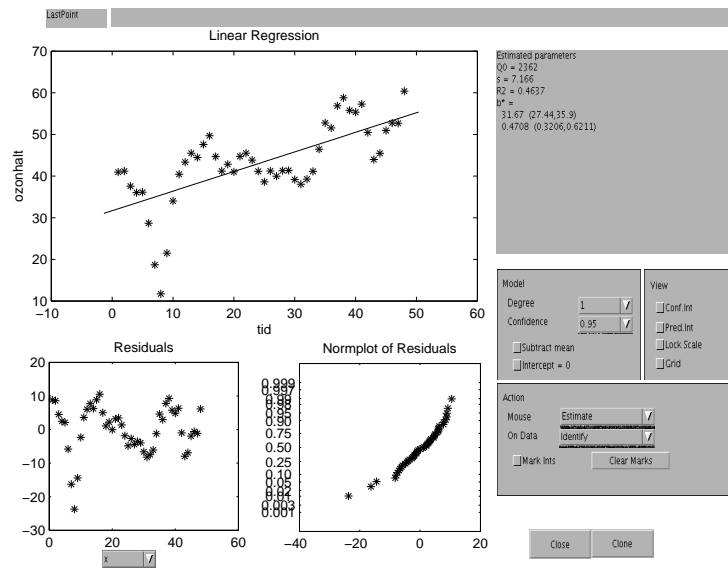
$$Y_t = \mu_t + R_t,$$

där μ_t är trenden och $\{R_t\}$ en stationär process, d.v.s. genomsnittsvärdet för $\{R_t\}$ är konstant under hela tidsperioden. Med denna uppdelning kan man t.ex. i μ_t få en ”jämn” (smooth) komponent som ska fånga upp säsongsmönster och/eller varierande trend och i $\{R_t\}$ en mer brusig komponent.

Om trenden μ_t tycks vara linjär kan den skattas med metoderna från enkel linjär regression, se avsnitt 2.3. Oftast är trenden emellertid en mera ”mjuk” funktion, då kan man använda ett glidande medelvärde. En utjämning (smoother) av $\{y_t\}$ fås t.ex. genom

$$\mu_t^* = \frac{y_{t-1} + y_t + y_{t+1}}{3}.$$

Den resulterande utjämnaren, μ_t^* , har skapats genom ett s.k. glidande medelvärde med fönster 3. Det betyder att man låter ett ”fönster” av tre mätningars bredd glida utmed tidsserien $\{y_t\}$ och att man för varje steg bildar medelvärdet av de observationer som syns i fönstret.



Figur 23: Linjär regressionsmodell anpassad till ozonhalterna. Observera utseendet på residualerna (nederst till vänster).

Exempel 5.3. Ozonhalt i Lund. Vi använder metoderna från regressionsavsnittet och anpassar en regressionslinje, $31.67 + 0.4708 \cdot \text{tid}$ (se figur 23) och gör därmed komponentuppdelningen ”linjär trend + brus”. Observera att den linjära modellens residualer (bruset), som plottas nederst till vänster, uppvisar ett tydligt mönster. I följande exempel beskrivs detta ytterligare. \square

Exempel 5.4. En utjämning med fönster 3 används på kiselhalterna, $\mu_t^* = \frac{y_{t-1} + y_t + y_{t+1}}{3}$ som sedan subtraheras från den ursprungliga tidsserien. I figur 24 visas komponentuppdelningen: överst markerar stjärnorna de ursprungliga mätningarna medan utjämnaren μ_t^* är markerad med streck. Underst visas det kvarvarande bruset efter subtraktionen. Observera att det är olika skalor på de två figurernas y-axlar. ”Trenden” μ_t består i stort sett av en säsongskomponent. \square

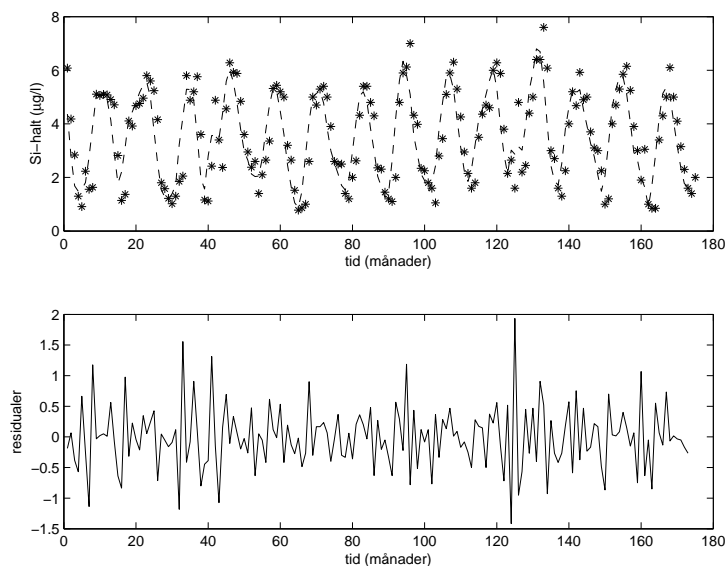
Idén ovan att dela upp processen $\{Y_t\}$ i en icke-stokastisk ”jämn” komponent μ_t och en stokastisk ”brusig” komponent $\{R_t\}$ är inte problemfri i praktiken. Det finns naturligtvis ingen entydig uppdelning. I exemplet med ozonhalterna kunde vi lika väl använt ett polynom av högre grad eller ett glidande medelvärde (eller något annat) för att skatta μ_t . Vad som väljs är ofta av en fråga om vad som ”räkar vara praktiskt”. Denna svaghet hindrar förstås inte att en uppdelning oftast ger värdefull information om den uppmätta tidsserien.

5.2.2 Beroende i tidsserien

I avsnitt 4.1 studerade vi samband mellan två stokastiska variabler X och Y och definierade kovarians och korrelationskoefficient mellan variablerna. Nu ska vi använda samma mått men på olika stokastiska variabler inom processen $\{Y_t\}$, vi tittar då på **autokorrelationen** i processen.

Antag att vi har en process $\{Y_t\}$ som är stationär, d.v.s. har en konstant genomsnittsnivå. Betrakta två slumpvariabler i processen, Y_{t-k} och Y_t som är på tidsavstånd k från varandra. De har kovarians $C(Y_{t-k}, Y_t)$ och korrelationskoefficient $\frac{C(Y_{t-k}, Y_t)}{C(Y_t, Y_t)} = \frac{C(Y_{t-k}, Y_t)}{V(Y_t)}$.

Om kovariansen (och därmed också korrelationskoefficienten) beror enbart på avståndet k och inte på tidpunkten t , sägs processen vara svagt stationär. I kiselhaltexemplet tänker vi oss alltså att när vi undersöker hur beroende en mätning är med en mätning fyra månader framåt, så är beroendet på ”fyramånadersavstånd” det samma oavsett när i tidsserien vi betraktar det.



Figur 24: Kiselhalter från Ljungbyån. Överst en utjämning (fönster 3), tillsammans med ursprungliga mätningar, underst det kvarvarande bruset.

Autokorrelationen på avstånd ("lag") k för en svagt stationär process är korrelationskoefficienten mellan mätningar som är på tidsavstånd k . Den är alltså ett mått på sambandet mellan variablerna Y_1 och Y_{1+k} samt mellan variablerna Y_2 och Y_{2+k} ..., samt mellan variablerna Y_{n-k} och Y_n . Alla dessa samband antas vara lika stora och korrelationskoefficienten betecknas ρ_k . Observera att ρ_0 alltid är 1. För övrigt gäller det som för vanliga korrelationskoefficienter att $-1 \leq \rho_k \leq 1$ och

positiv samvariation (positiv korrelation) på tidsavstånd (lag) k	motsvarar	$\rho_k > 0$
negativ samvariation (negativ korrelation) på tidsavstånd (lag) k	motsvarar	$\rho_k < 0$
ingen samvariation (ingen korrelation) på tidsavstånd (lag) k	motsvarar	$\rho_k = 0$

En kommentar om beteckningar. I avsnitt 4.1 betecknade vi korrelationskoefficienten ρ_{xy} eftersom det uttryckte sambandet mellan talparen $(X_1, Y_1), (X_2, Y_2), (X_3, Y_3), \dots, (X_n, Y_n)$. Nu tittar vi på samband inom tidsserien $\{Y_t\}$. För att t.ex. studera sambandet på tidsavstånd 3 ska vi betrakta talparen $(Y_1, Y_4), (Y_2, Y_5), \dots, (Y_{n-3}, Y_n)$ och betecknar sambandet ρ_3 .

5.2.3 Skattning av autokorrelationsfunktionen

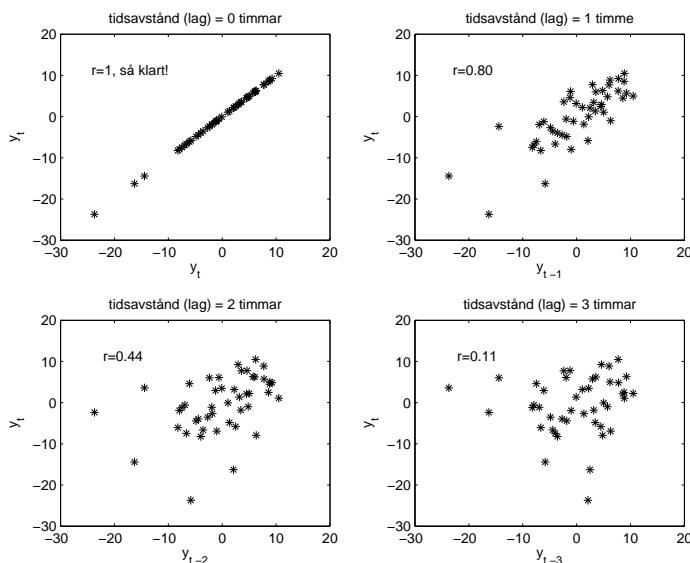
När ρ_k ska skattas kan man alltså beräkna korrelationskoefficienten r_k för de $n - k$ talparen $(y_1, y_{1+k}), \dots, (y_t, y_{t+k}), \dots, (y_{n-k}, y_n)$. I analogi med skattningen r_{xy} från avsnitt 4.1 får vi

$$r_k = \frac{c_k}{s_y s_y} = \frac{\sum_{t=1}^{n-k} (y_t - \bar{y})(y_{t+k} - \bar{y})}{\sqrt{\sum_{t=1}^n (y_t - \bar{y})^2} \sqrt{\sum_{t=1}^n (y_t - \bar{y})^2}} = \frac{\sum_{t=1}^{n-k} (y_t - \bar{y})(y_{t+k} - \bar{y})}{\sum_{t=1}^n (y_t - \bar{y})^2}.$$

Observera att i uttrycket för r_k består täljaren av en summa av $n - k$ termer medan nämnarens summa har n . Oftast gör man emellertid ingen korrigerig för att ta hänsyn till detta, det olika antalet termer spelar ingen större roll om k är tillräckligt liten i förhållande till n , vilket det i praktiken alltid är.

Genom att skatta ρ_k för olika värden på tidsavstånden k och plotta dessa skattningar mot k har man plottat den s.k. **skattade autokorrelationsfunktionen**. Med den får man en samlad bild av korrelationskoefficienterna på olika tidsavstånd.

Exempel 5.5. Residualerna när vi anpassat en linjär regressionsmodell till ozonmätningarna visas nederst till vänster i figur 23. De visar ett tydligt mönster som tyder på ett beroende



Figur 25: Residualer från ozonmätningarna, residualer på olika tidsavstånd (lag) plottade mot varandra.

mellan näraliggande mätningar. I figur 25 har vi plottat talpar på avstånd k , för $k = 0, 1, 2, 3$ och motsvarande skattning av korrelationskoefficient är markerad i figurerna. I plotten längst ner till höger, där vi studerar beroende på tidsavstånd (lag) 3 hos residualerna $\{y_t\}$, är alltså de 45 talparen $(y_1, y_4), (y_2, y_5), \dots, (y_{45}, y_{48})$ utritade.

Mer överskådligt är det i figur 26 där korelationskoefficienterna sammanförs i den skattade autokorrelationsfunktionen. Sammantaget verkar det finnas ett starkt positivt samband mellan ozonmätningar som är gjorda omedelbart efter varandra, på en timmas tidsavstånd. Sambandet är lite svagare mellan mätningar på två timmars avstånd för att avklinga på längre tidsavstånd. Här har vi valt att enbart titta på ett fåtal k -värden eftersom tidsserien är förhållandevis kort.

□

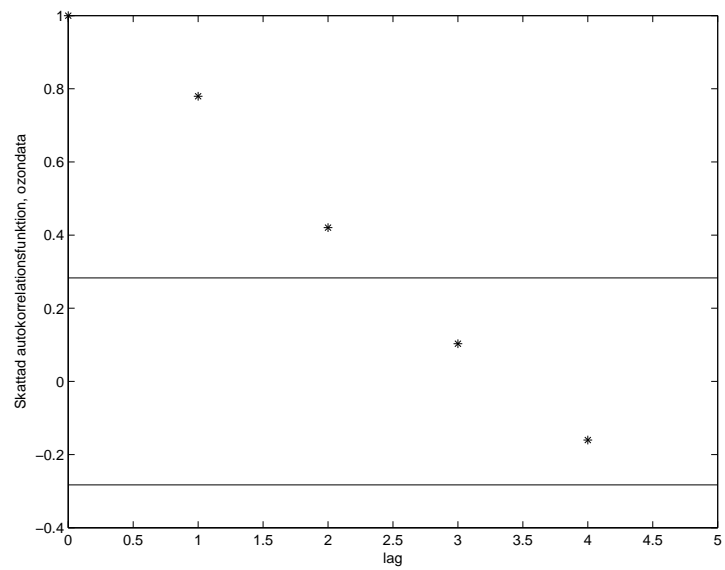
Hur mycket ska r_k avvika från 0 för att man kan anse att det finns ett samband mellan mätningar på avstånd k ? Svaret borde bl.a. bero på tidsseriens längd, d.v.s. n . Ett första grovt test för att avgöra om autokorrelationerna är signifikant skilda från 0 är följande: Om Y_1, \dots, Y_n är oberoende (vitt brus) gäller att $\rho_k^* \lesssim N(0, \frac{1}{n})$ dvs om $\rho_k = 0$ så ska endast 5% av ρ_k^* hamna utanför gränserna $(0 - \lambda_{0.025} \frac{1}{\sqrt{n}}, 0 + \lambda_{0.025} \frac{1}{\sqrt{n}})$. Dessa två ”kontrollgränser” ritas därför ofta ut i plotten över den skattade autokorrelationsfunktionen. Om alltså ρ_k^* hamnar utanför gränserna bör vi betvivla att $\rho_k = 0$. Använd testet med en viss försiktighet eftersom de olika skattningarna ρ_k^* är beroende.

Exempel 5.6. I figur 27 skattas autokorrelationsfunktion på ursprungliga kiselhalterna (överst) samt på bruset i komponentuppdelningen (nederst). Medan ursprungliga Si-halter uppvisar en tydlig säsongsmönster verkar det inte finnas något kvar i bruset, några r_k (ca 5%) ”får ligga utanför gränserna”. Observera emellertid att den skattade autokorrelationen på lag 1 är starkt negativ. Detta beror på att vi som utjämning använt ett glidande medelvärde. □

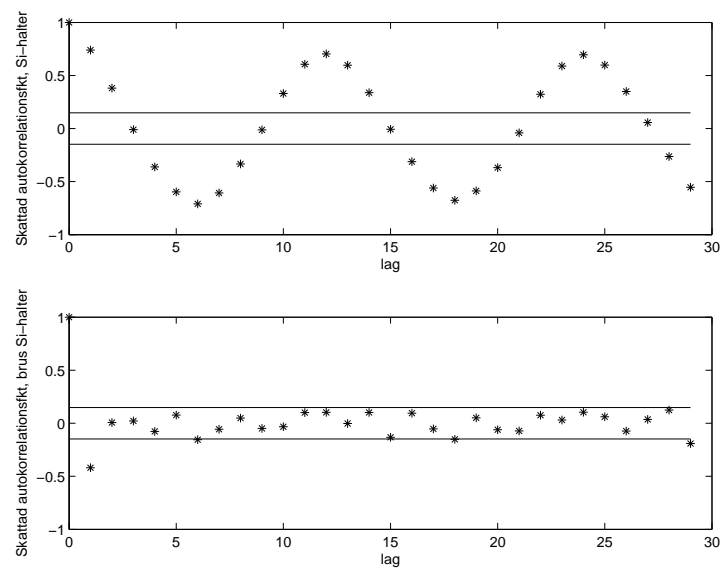
5.2.4 Matlabkommandon för skattning av autokorrelationsfunktionen

För beräkningarna av autokorrelationsfunktionen ska ni utnyttja kommandot `corr` i Matlab.

- Det vi vill beräkna är uttrycket $\rho_k^* = r_k = \frac{\sum (y_t - \bar{y})(y_{t+k} - \bar{y})}{\sum (y_t - \bar{y})^2}$ för olika värden på k .



Figur 26: Skattad autokorrelationsfunktion för residualerna i ozondata då en linjär regressionsmodell anpassats.



Figur 27: Skattad autokorrelationsfunktion på ursprungliga kiselhalterna (överst) samt på bruset i komponentuppdelningen (nederst).

- Om data finns i datavektorn y ger kommandot $c = \text{corr}(y - \text{mean}(y), M)$ en vektor c bestående av M element vilket ger den sökta skattade autokorrelationsfunktionen ρ_k^* för värdena $k = 0, 1, \dots, M - 1$.
- För att få autokorrelationsfunktionen utritad använder ni lämpligen plotkommandot $\text{plot}([0:M-1], c, '*)$.

Exempel 5.7. Antag att kiselhalterna ligger i vektorn si . För att rita ut den skattade autokorrelationsfunktionen överst i figur 27 använde vi de två Matlabkommandona $c = \text{corr}(si - \text{mean}(si), 30); \text{plot}([0:29], c, '*)$ □

5.3 Modeller

Med hjälp av en skattad autokorrelationsfunktion kan vi beskriva hur beroendet ser ut i en tidsserie. Nästa steg är att försöka göra en matematisk modell för beroendet. En AR(1)-process (first-order autoregressive process) är en enkel slumpmodell som ofta används vid modellering av miljötidsserier.

5.3.1 AR(1)-processer

Om tidsserien $\{y_t\}$ är en realisering (observation) av processen $\{Y_t\}$ tänker vi oss att mätningen vid tidpunkt t , Y_t kan beskrivas enligt

$$Y_t = \alpha \cdot Y_{t-1} + e_t$$

där $-1 < \alpha < 1$ och slumpvariablerna $\{e_t\}$ är oberoende med varians σ^2 .

Faktorn α benämnes ibland som "minnesfaktor" eftersom den talar om hur mycket av mätningen från tidpunkt $t - 1$ man "kommer ihåg" då man gör mätningen vid tidpunkt t . Till detta adderas sedan en "färsk slump", e_t , med väntevärde 0 och standardavvikelse σ . I simuleringar i detta kompendium har vi för enkelhets skull oftast antagit att $\{e_t\}$ är oberoende och $N(0, 1)$.

Exempel 5.8. I en AR(1)-process med $\alpha = 0$ gäller $Y_t = e_t$, där e_t är oberoende. Observationerna $\{y_t\}$ är alltså helt oberoende - man säger att man har "vitt brus". Figur 28 visar en realisering av en sådan process där $\{e_t\} \in N(0, 1)$ samt motsvarande skattade autokorrelationsfunktion. Eftersom vi har oberoende observationer är $\rho_k = 0$, $k = 1, 2, \dots$ och motsvarande skattningar ρ_k^* ligger också nära 0. □

För en AR(1)-process kan man visa att autokorrelationsfunktionen beskrivs av ett enkelt samband,

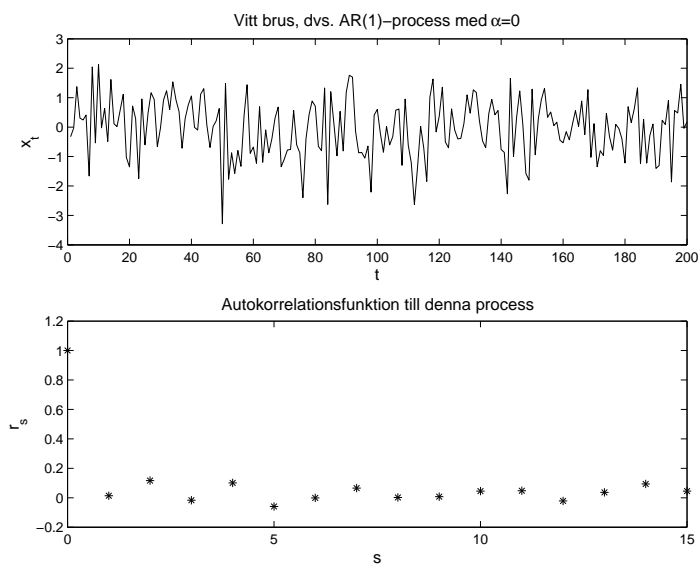
$$\rho_k = \alpha^k, \quad k = 0, 1, 2, \dots$$

där $-1 < \alpha < 1$.

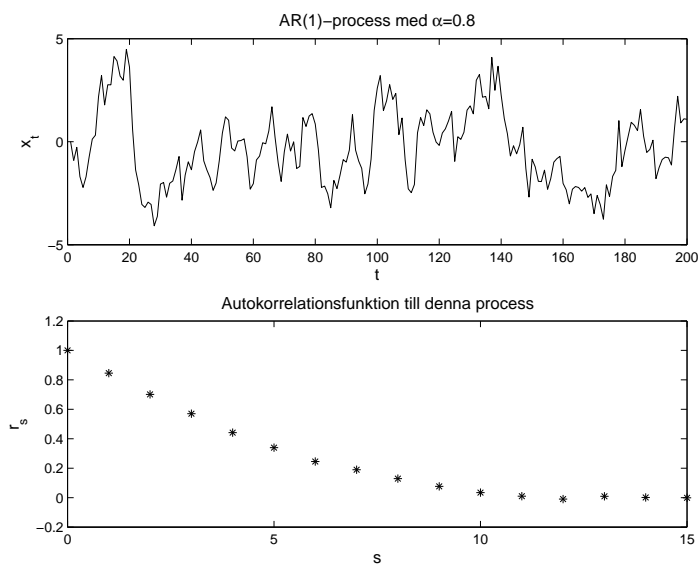
Exempel 5.9. Med en minnesfaktor på $\alpha = 0.8$ har vi modellen $Y_t = 0.8 \cdot Y_{t-1} + e_t$, där e_t är oberoende. För denna process gäller $\rho_k = 0.8^k$, $k = 0, 1, 2, \dots$, d.v.s. sambandet mellan näraliggande mätningar är högt (0.8 på tidsavstånd 1, 0.64 på tidsavstånd 2 o.s.v.) för att sakta avklinga. Figur 29 visar en realisering av en sådan process med $e_t \in N(0, 1)$, samt motsvarande skattade autokorrelationsfunktion. □

Observera "trögheten" i tidsserien i figur 29. AR(1)-processer med positiva minnesfaktorer är ofta en lämplig modell för att modellera tidsserier med en inneboende "biologisk tröghet", ju högre värde på minnesfaktorn desto "trögare" system. I det tidigare exemplet med ozonhalt verkar det inte otroligt att residualerna skulle kunna beskrivas med en AR(1)-process med positiv minnesfaktor, jfr figur 26.

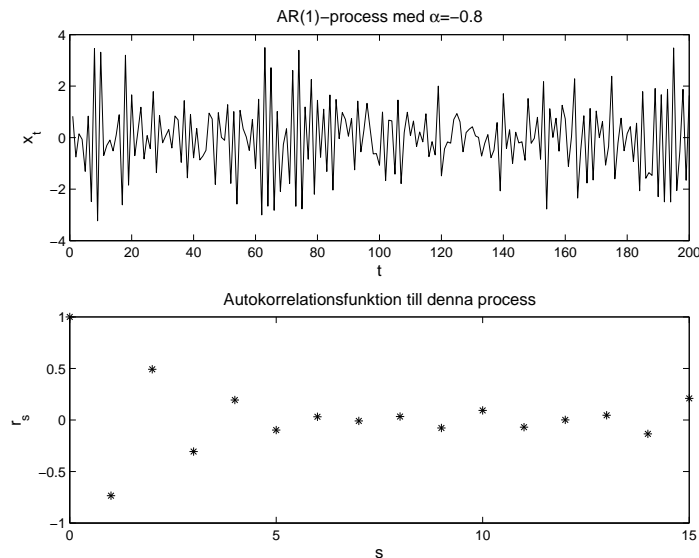
Exempel 5.10. I en sjö med ett tillflöde mäts halten av en viss förorening en gång i månaden. Det finns inga utsläpp i sjön, utan föroreningarna kommer från tillflödet. Under en månad byts 30% av vattnet i sjön ut. Låt $\{X_t\}$ vara halten i sjön vid tidpunkt t och $\{Z_t\}$ vara den



Figur 28: AR(1)-process med $\alpha = 0$, vitt brus.



Figur 29: AR(1)-process med $\alpha = 0.8$.



Figur 30: AR(1)-process med $\alpha = -0.8$.

genomsnittliga halten av inflödet under den senaste månaden, d.v.s. från $t - 1$ till t . Då gäller, så när som på mätfel, $X_t = 0.7 \cdot X_{t-1} + 0.3 \cdot Z_t$.

Antag att $\{Z_t\}$ är oberoende stokastiska variabler med väntevärde μ och varians σ_z^2 . Inför de nya variablerna $Y_t = X_t - \mu$ och $e_t = 0.3 \cdot (Z_t - \mu)$. Då gäller det att $\{e_t\}$ är oberoende med $E(e_t) = 0$ och $V(e_t) = 0.3^2 \cdot \sigma_z^2$, d.v.s. $\{e_t\}$ är vitt brus. Dessutom är $Y_t = X_t - \mu = 0.7 \cdot (X_{t-1} - \mu) + 0.3 \cdot (Z_t - \mu) = 0.7 \cdot Y_{t-1} + e_t$.

Sammantaget har vi visat att $\{Y_t\}$ kan beskrivas som en AR(1)-process med $\alpha = 0.7$. \square

Exempel 5.11. Med en minnesfaktor på $\alpha = -0.8$ i en AR(1)-process har vi modellen $Y_t = -0.8 \cdot Y_{t-1} + e_t$, där e_t är oberoende. För denna process gäller $\rho_k = (-0.8)^k$, $k = 0, 1, 2, \dots$, d.v.s. sambandet mellan näraliggande mätningar alternerar mellan positiva och negativa korrelationer. Figur 30 visar en realisering av en sådan process med $e_t \in N(0, 1)$, samt motsvarande skattade autokorrelationsfunktion. Observera hur snabbt tidsserien oscillerar kring nollnivån. Det är svårare att hitta praktiska "miljötilämpningar" som beskrivs av AR(1)-processer med negativa minnesfaktorer. \square

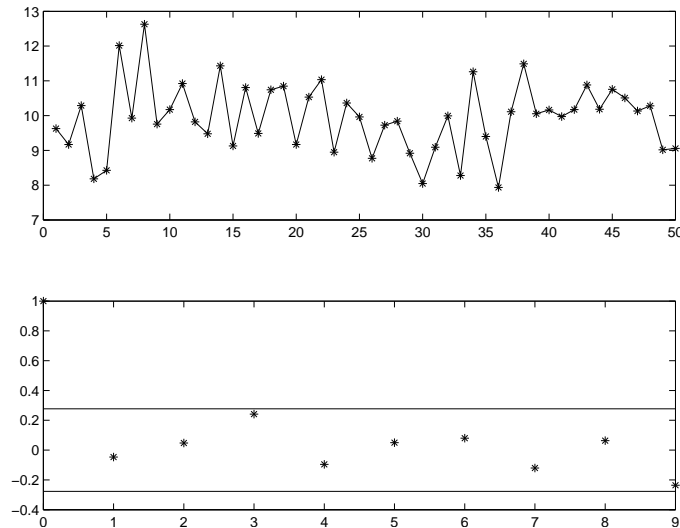
5.3.2 Simulering av AR(1)-processer i Matlab

Om man vill simulera mätningar från en AR(1)-process i Matlab utnyttjas kommandot `filter`. Antag att vi vill ha n observationer från processen $Y_t = \alpha \cdot Y_{t-1} + e_t$ där $e_t \in N(0, 1)$. Genom kommandot `e=normrnd(0,1,1,n)` skapas de n "färska slumpalen" $\{e_t\}$. Matlabkommandot `y=filter(1, [1, -alpha], e)` ger sedan den önskade tidsserien. Observera det negativa tecknet framför värdet på α .

Den första 1 i `filter` hör ihop med att det är just en AR(1)-process (man kan tänka sig ett mer komplicerat beroende) medan uttrycket `[1, -alpha]` används eftersom det är vänstra ledets koefficienter i relationen $y_t - \alpha \cdot y_{t-1} = e_t$.

5.4 Beroende mätningar påverkar analysen

Antag att vi har n oberoende mätningar y_1, \dots, y_n , där motsvarande s.v. alla har väntevärde μ och varians σ^2 . En skattning av μ är $\mu_{ober}^* = \bar{y}$, där indexet *ober* står för att observationerna är oberoende. Variansen av denna



Figur 31: Vitt brus adderat till konstanten 10.

skattning är $V(\mu_{ober}^*) = \frac{\sigma^2}{n}$.

Antag nu att mätningarna y_1, \dots, y_n är beroende så att på tidsavstånd k är autokorrelationen ρ_k . Vi skattar fortfarande μ med medelvärdet, $\mu_{ber}^* = \bar{y}$, men mätningarna är nu beroende. Då kan man visa att $V(\mu_{ber}^*) = \frac{\sigma^2}{n} + \frac{2\sigma^2}{n^2} \sum_{k=1}^{n-1} (n-k)\rho_k$. Eftersom man ofta har ett positivt beroende i "miljötidsserier" innebär det att $\rho_k > 0$ och därmed $V(\mu_{ber}^*) > V(\mu_{ober}^*)$. Det är alltså direkt fel att behandla beroende mätningar som om de var oberoende. Relationen ovan säger t.ex. att om vi "slentrianmässigt" använder formeln $\frac{\sigma^2}{n}$, baserad på de beroende mätningarna, för att skatta $V(\mu_{ber}^*)$ blir skattningen för liten.

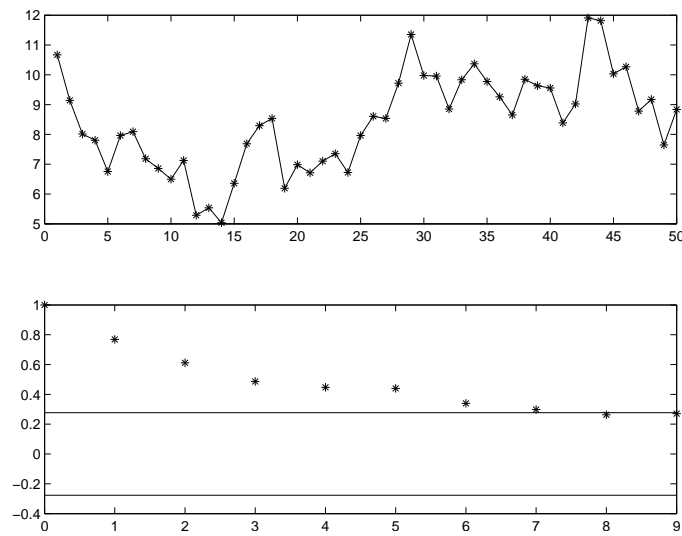
Ett annat sätt att uttrycka det är att n beroende observationer ger inte så mycket information som n oberoende observationer. Misstänker man att man har beroende i sin tidsserie, eller har visat det i sin skattade autokorrelationsfunktion, ska man alltså inte göra för täta mätningar. Det är bättre att "sampla" glesare i tiden så att man får oberoende, eller nästan oberoende, observationer. Alternativet är att använda sig av komplicerade formler för beroende data.

5.4.1 Beroende data påverkar trendanalysen!

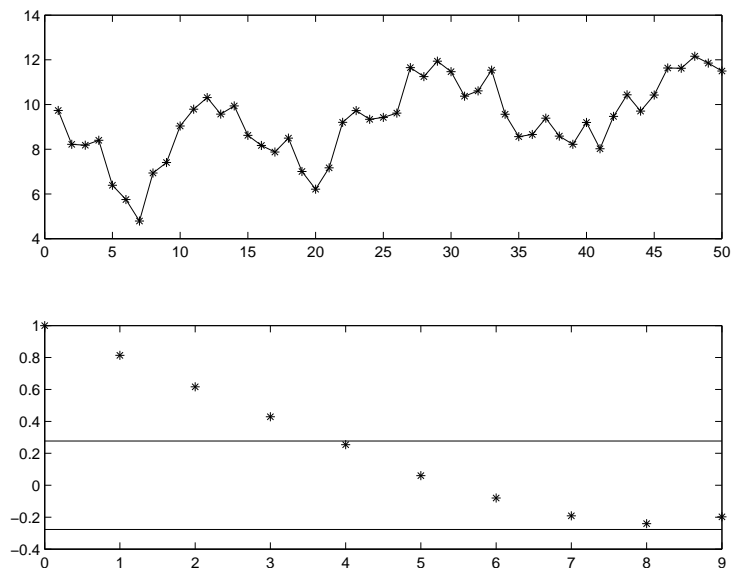
Även vid trendanalys ger beroende data problem. En ökad variabilitet med ökat t kan få tidsserien att "driva" och ge ett falskt intryck av trend eller mönster.

Exempel 5.12. I figur 31 är en "miljötidsserie", $\{y_t\}$, simulerad, vilken består av 50 oberoende slumpstal $N(0, 1)$ adderade till konstanten 10. "Mätningarna" är alltså vitt brus kring konstanten 10 och ingen förändring i trend har skett under perioden. Den skattade autokorrelationsfunktionen indikerar också att mätningarna är oberoende.

Antag nu att vi i stället till konstanten 10 adderar tidsserien, $\{a_t\}$, på 50 mätningar som består av beroende mätfel, där beroendet är av typen AR(1) med $\alpha = 0.8$. För vår "uppmätta" tidsserie $\{y_t\}$ gäller då $y_t = 10 + a_t$ där $a_t = 0.8 \cdot a_{t-1} + e_t$ med $\{e_t\}$ oberoende $N(0, 1)$. Det har alltså fortfarande inte hänt något med "miljövariabeln", den är fortfarande konstant 10, men mätfelen är beroende. I figur 32 och figur 33 ses två olika realiseringar av $\{y_t\}$. Beroende på vad Matlab råkat ge oss för slumpstal kan tidsserierna bli av lite olika art. I den första realiseringen syns en mjukt varierande trend och i den andra kan man skönja ett cykliskt mönster. Inget av dessa "trender" eller "cykler" finns där - det är det starka beroendet mellan mätfelen eller annorlunda uttryckt den stora "trögheten" i miljötidsserien som kan ge upphov till "falsa" mönster.



Figur 32: En realisering, med tillhörande skattad autokorrelationsfunktion, av en AR(1)-process med $\alpha = 0.8$ adderad till konstanten 10.



Figur 33: Ytterligare en realisering, med tillhörande skattad autokorrelationsfunktion, av en AR(1)-process med $\alpha = 0.8$ adderad till konstanten 10.

5.5 Läs mer om trendanalys och tidsserier

För den som vill läsa mer om trendanalys i miljötidsserier rekommenderas följande nyutkomna bok:

Chandler, R.E and Scott M: *Statistical Methods for Trend Detection and Analysis in the Environmental Sciences*, Wiley, 2011.

Den finns bl.a. tillgänglig i Matematikbiblioteket i Lund.

□

6 Mer om trendanalys

I tidigare avsnitt har trendanalysen byggts på en regressionsmodell $y_t = \alpha + \beta \cdot t + e_t$, där slumpfelen e_1, \dots, e_t är oberoende och normalfördelade.

Alternativt har vi genom en utjämning, t.ex. ett glidande medelvärde, fått ett beskrivande mått på trenden. Denna metod är speciellt användbart för beroende observationer men kräver oftast att observationerna är mätta regelbundet.

Det är emellertid inte ovanligt att miljötidsserier uppvisar en rad egenskaper som gör att ovanstående metoder inte är så lämpliga. Mätningarna kan t.ex.

- bestå enbart av ett fåtal värden;
- innehålla många saknade värden;
- innehålla outliers;
- innehålla mätningar under en detektionsnivå;
- bestå av värden där en normalfördelning, eller någon annan standardfördelning, ej är rimlig att anta.

Ett test för trend som kan hantera dessa problem är **Mann-Kendalls test**, som ibland används som ett första ”snabbtest” för att undersöka trend. Av någon anledning har den blivit speciellt populär bland forskare som studerar vattenkvalitetsdata men metoderna är naturligtvis generella att kunna användas på vilken sorts miljödata som helst. Det avgörande för testets popularitet är att det är s.k. fördelningsfritt (ibland säger man icke-parametriskt), dvs det förutsätter inget om den bakomliggande fördelningen hos data. Däremot förutsätter testet, i den form som det beskrivs här, att data är oberoende. För beroende observationer finns varianter på testet.

6.1 Mann-Kendalls test

Den uppmätta tidsserien betecknas som tidigare y_1, y_2, \dots, y_n .

Tanken är att bilda alla tänkbara par av observationer (y_j, y_k) där $k > j$ (det finns $\binom{n}{2}$ st) och jämföra antalet par där $y_k > y_j$ (positiv trend) med antalet par där $y_k < y_j$ (negativ trend).

Som testkvantitet S tas då

$$S = (\text{antalet par } (y_j, y_k) \text{ där } y_k > y_j) - (\text{antalet par } (y_j, y_k) \text{ där } y_k < y_j).$$

I formler kan detta uttryckas med hjälp av sgn-funktionen. Låt

$$\text{sgn}(y_k - y_j) = \begin{cases} 1, & y_k - y_j > 0; \\ 0, & y_k - y_j = 0; \\ -1, & y_k - y_j < 0. \end{cases}$$

Den aktuella testkvantiteten S kan då uttryckas som

$$S = \sum_{j=1}^{n-1} \sum_{k=j+1}^n \text{sgn}(y_k - y_j).$$

Exempel 6.1. Tidsserien består av de 6 observationerna

5 4 7 9 8 12

Antal par som ska jämföras är $\binom{6}{2} = 15$, nämligen

(5,4)	(5,7)	(5,9)	(5,8)	(5,12)
	(4,7)	(4,9)	(4,8)	(4,12)
		(7,9)	(7,8)	(7,12)
			(9,8)	(9,12)
				(8,12)

Av dessa par är det 13 där andrakoordinaten är större än förstakoordinaten och 2 med det omvända förhållandet. Detta ger

$$S = 13 - 2 = 11.$$

□

Observera att för en tidsserie med n observationer kan S högst bli $\binom{n}{2}$ vilket sker om det för alla observationer gäller att en observation är större än den omedelbart föregående, dvs serien har en stigande trend hela tiden. På samma sätt kan S inte understiga $-\binom{n}{2}$, vilket motsvarar en strängt avtagande trend. Om däremot ingen trend finns i serien tenderar S att vara kring 0. Nollhypotesen H_0 : ”ingen trend” bör alltså förkastas om S avviker allt för mycket från 0. Med hjälp av tabell (finns i slutet av detta avsnitt) kan testet göras med hjälp av direktmetoden som följande exempel visar.

Exempel 6.2. I exemplet ovan blev $S = 11$ för en tidsserie på $n = 6$. Antag att vi misstänker att det finns en positiv trend i data, vår mothypotes är alltså ensidig. Tabellen ger att under förutsättning att ingen trend finns är $P(S \geq 11) = \alpha_0 = 0.028$. Tolkningen av detta α_0 (kallas ibland för P-värde) är att hypotesen om ”ingen trend” kan förkastas på nivå 0.028. Vi har alltså, på signifikansnivå 0.028, påvisat en positiv trend i data. □

För stora värden på n , ($n > 10$), kan man använda att då H_0 är sann gäller att S är approximativt normalfördelad med väntevärde 0 och varians $\frac{1}{18}n(n-1)(2n+5)$. Det gäller alltså att

$$S \in N \left(0, \sqrt{\frac{1}{18}n(n-1)(2n+5)} \right).$$

Exempel 6.3. I en tidsserie om $n = 12$ värden misstänker man en negativ trend i data. Man har beräknat ett värde på testkvantiteten, $S = -23$. Är detta förenligt med en normalfördelning med väntevärde 0 och varians $\frac{1}{18}12(12-1)(2 \cdot 12+5) = 212.667$, dvs med $N(0, \sqrt{212.667})$?

Vi undersöker genom att beräkna $t = \frac{-23-0}{\sqrt{212.667}} = -1.58$ och konstaterar att eftersom $-1.58 > -\lambda_{0.05} = -1.65$ kan H_0 ej förkastas på nivå 0.05.

Alternativt använder vi direktmetoden och beräknar $\alpha_0 = P(S < -23 \text{ om } S \in N(0, 212.667)) = 0.0574$. Det senare beräknas enklast via kommandot `normcdf(-23, 0, sqrt(12*11*29/18))` i Matlab. Eftersom P-värdet överstiger ”standardnivån” 0.05 är slutsatsen som tidigare att nollhypotesen ej kan förkastas på nivå 0.05. Om vi påstår att nollhypotesen ska förkastas, dvs att det finns en negativ trend, har vi alltså en felrisk på 0.0574. □

6.2 Skattning av trenden

Om man konstaterat att en trend föreligger vill man förmodligen skatta trenden också. En sådan skattning är Sens skattning av trenden β . För varje par (y_j, y_k) där $k > j$ bildas kvoten

$$b_{jk} = \frac{y_k - y_j}{k - j}$$

som kan betraktas som "riktningskoefficienten" mellan tidpunkterna j och k . Sens skattning är medianen av alla dessa $\binom{n}{2}$ skattningar:

$$\beta^* = \text{median}\left(\frac{y_k - y_j}{k - j}\right).$$

Exempel 6.4. Om tidsserien ovan är uppmätt vid tidpunkterna 1, 2, 3, 5, 6, 7 (vid tidpunkt 4 kunde ingen mätning göras) har vi

tidpunkt:	1	2	3	5	6	7
y_t :	5	4	7	9	8	12

och de 15 b_{jk} -värdena är

$$\begin{array}{llllll} \frac{(4-5)}{(2-1)} = -1 & \frac{(7-5)}{(3-1)} = 1 & \frac{(9-5)}{(5-1)} = 1 & \frac{(8-5)}{(6-1)} = 0.6 & \frac{(12-5)}{(7-1)} = 1.17 \\ & \frac{(7-4)}{(3-2)} = 3 & \frac{(9-4)}{(5-2)} = 1.67 & \frac{(8-4)}{(6-2)} = 1.33 & \frac{(12-4)}{(7-2)} = 1.6 \\ & & \frac{(9-7)}{(5-3)} = 1 & \frac{(8-7)}{(6-3)} = 0.33 & \frac{(12-7)}{(7-3)} = 1.25 \\ & & & \frac{(8-9)}{(6-5)} = -1 & \frac{(12-9)}{(7-5)} = 1.5 \\ & & & & \frac{(12-8)}{(7-6)} = 4 \end{array}$$

Medianen av dessa 15 värden är 1.17 så skattningen av trenden är $\beta^* = 1.17$. \square

6.3 Seasonal Kendall test

Om tidsserien uppvisar säsongvariation (eller annan cyklisk variation) kan man inte använda metoderna rakt av. Ett trick är då att betrakta varje säsong för sig, beräkna ett S -värde, S_i , för varje säsong, och sedan summera S_i -värdena till en total testkvantitet, S_{total} . Den totala testkvantiteten

$$S_{total} = \sum S_i$$

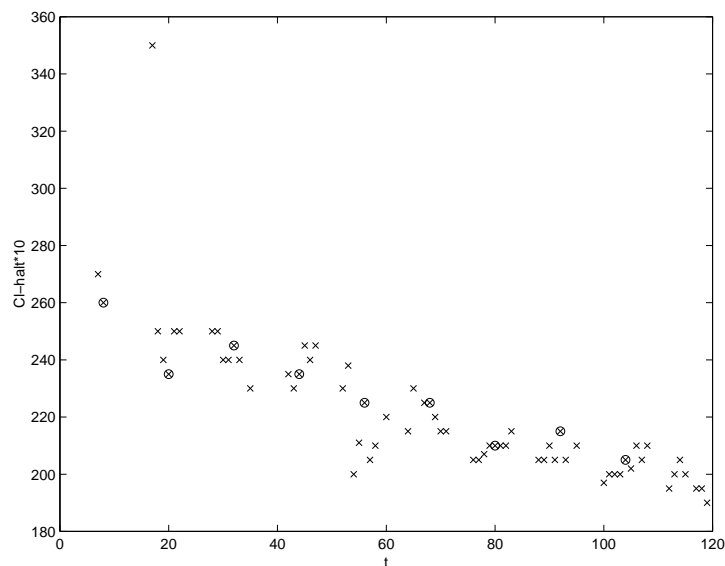
är approximativt normalfördelad $N\left(0, \sqrt{\sum_i \frac{1}{18} n_i(n_i - 1)(2n_i + 5)}\right)$ där n_i är antalet observationer under säsong i .

Exempel 6.5. Vid fem olika stationer på den norra sidan av Lake Erie, Ontario i Kanada mäts bl.a. kloridhalten (mg/l) under en längre period. Vid station 501 "Long Point Bay" fick man följande resultat (medianvärdet multiplicerat med 10).

År	jan	feb	mar	apr	maj	jun	jul	aug	sep	okt	nov	dec
1970							270	260				
1971					350	250	240	235	250	250		
1972				250	250	240	240	245	240		230	
1973						235	230	235	245	240	245	
1974				230	238	200	211	225	205	210		220
1975				215	230		225	225	220	215	215	
1976				205	205	207	210	210	210	210	215	
1977				205	205	210	205	215	205		210	
1978				197	200	200	200	205	202	210	205	210
1979				195	200	205	200		195	195	190	

I figur 34 är tidsserien utritad, augusti månads mätningar är markerade med en ring runt krysen.

Tanken är att beräkna ett S_i för var och en av de nio månaderna. För augusti månad gäller att av de totalt $\binom{9}{2} = 36$ talparen som ska studeras har 2 par en andrakoordinat som överstiger förstakoordinaten, 32 par har det omvända förhållandet medan i 2 par är de två koordinaterna



Figur 34: Kloridhalt vid station "Long Point Bay" under perioden 1972-1979, augusti månads mätningar är markerade.

lika. För denna månad gäller alltså att $S_i = 2 - 32 = -30$. Antalet observationer denna månad är $n_i = 9$ vilket ger variansen $\frac{1}{18}n_i(n_i - 1)(2n_i + 5) = \frac{1}{18}9(9 - 1)(2 \cdot 9 + 5) = 92$.

Motsvarande beräkningar görs för de övriga 8 månaderna, vilket är sammanställt i följande tabell.

månad	n_i	S_i	$V(S_i) = \frac{1}{18}n_i(n_i - 1)(2n_i + 5)$
1 (april)	7	$(0 - 20) = -20$	44.33
2 (maj)	8	$(0 - 26) = -26$	65.33
3 (juni)	8	$(5 - 22) = -17$	65.33
4 (juli)	10	$(1 - 42) = -41$	125
5 (aug)	9	$(2 - 32) = -30$	92
6 (sep)	9	$(3 - 32) = -29$	92
7 (okt)	7	$(1 - 17) = -16$	44.33
8 (nov)	7	$(1 - 19) = -18$	44.33
9 (dec)	2	$(0 - 1) = -1$	1
		$S_{total} = -198$	$V(S_{total}) = 573.65$

För att testa hypotesen "ingen trend" beräknas testkvantiteten $\frac{S_{total} - 0}{\sqrt{V(S_{total})}} = \frac{-198}{\sqrt{573.65}} = -8.27$. Eftersom absolutbeloppet av testkvantiteten överstiger $z_{0.0005} = 3.29$ förkastas hypotesen på signifikansnivå 0.001.

□

Tabell över sannolikheter för Mann-Kendalls test för trend

S	n=4	n=5	n=8	n=9	S	n=6	n=7	n=10
0	0.625	0.592	0.548	0.540	1	0.500	0.500	0.500
2	0.375	0.408	0.452	0.460	3	0.360	0.386	0.431
4	0.167	0.242	0.360	0.381	5	0.235	0.281	0.364
6	0.042	0.117	0.274	0.306	7	0.136	0.191	0.300
8		0.042	0.199	0.238	9	0.068	0.119	0.242
10		0.0 ² 83	0.138	0.179	11	0.028	0.068	0.190
12			0.089	0.130	13	0.0 ² 83	0.035	0.146
14			0.054	0.090	15	0.0 ² 14	0.015	0.108
16			0.031	0.060	17		0.0 ² 54	0.078
18			0.016	0.038	19		0.0 ² 14	0.054
20			0.0 ² 71	0.022	21		0.0 ³ 20	0.036
22			0.0 ² 28	0.012	23			0.023
24			0.0 ³ 87	0.0 ² 63	25			0.014
26			0.0 ³ 19	0.0 ² 29	27			0.0 ² 83
28			0.0 ⁴ 25	0.0 ² 12	29			0.0 ² 46
30				0.0 ³ 43	31			0.0 ² 23
32				0.0 ³ 12	33			0.0 ² 11
34				0.0 ⁴ 25	35			0.0 ³ 47
36				0.0 ⁵ 28	37			0.0 ³ 18
					39			0.0 ⁴ 58
					41			0.0 ⁴ 15
					43			0.0 ⁵ 28
					45			0.0 ⁶ 28

0.0²83 står för 0.0083

Exempel 6.6. I en tidsserie om 8 mätningar gäller - under förutsättning av ingen trend finns i data - $P(S \geq 14) = 0.054$. På grund av symmetri gäller, under samma förutsättning, även att $P(S \leq -14) = 0.054$.

□

7 Appendix: ML- och MK skattningar av parametrarna i enkel linjär regression

7.1 Några hjälpresultat

Vi börjar med ett par användbara beteckningar och räkneregler för de summor och kvadratsummor som kommer att ingå i skattningarna. Då alla summor nedan löper från 1 till n avstår vi från att skriva ut summationsindexen.

Först har vi att en ren summa av avvikelser av ett antal observationer kring sitt medelvärde är noll

$$\sum (x_i - \bar{x}) = \sum x_i - n\bar{x} = [\bar{x} = \frac{1}{n} \sum x_i] = \sum x_i - \sum x_i = 0 \quad (1)$$

Några beteckningar för kvadratiska- och korsavvikelser kring medelvärde

$$S_{xx} = \sum (x_i - \bar{x})^2, \quad S_{xy} = \sum (x_i - \bar{x})(y_i - \bar{y}), \quad S_{yy} = \sum (y_i - \bar{y})^2$$

där vi känner igen den första och sista från stickprovsvarianserna för x resp. y , $s_x^2 = S_{xx}/(n-1)$ och motsvarande för y . Dessa summor kan skrivas på ett antal former, t.ex kan S_{xy} utvecklas till

$$\begin{aligned} S_{xy} &= \sum (x_i - \bar{x})(y_i - \bar{y}) = \sum x_i(y_i - \bar{y}) - \bar{x} \sum (y_i - \bar{y}) = \sum x_i(y_i - \bar{y}) \quad \text{eller} \\ S_{xy} &= \sum (x_i - \bar{x})(y_i - \bar{y}) = \sum (x_i - \bar{x})y_i - \bar{y} \sum (x_i - \bar{x}) = \sum (x_i - \bar{x})y_i \end{aligned}$$

där sista summan i andra leden blir noll enligt (1). Motsvarande räkneregler gäller för S_{xx} och S_{yy} och vi har sammanfattningsvis

$$S_{xy} = \sum (x_i - \bar{x})(y_i - \bar{y}) = \sum x_i(y_i - \bar{y}) = \sum (x_i - \bar{x})y_i \quad (2)$$

$$S_{xx} = \sum (x_i - \bar{x})^2 = \sum x_i(x_i - \bar{x}) \quad \text{och motsvarande för } S_{yy} \quad (3)$$

7.2 Punktskattningar

ML-skattning av α , β och σ^2 då y_i är oberoende observationer av $Y_i \in N(\alpha + \beta x_i, \sigma)$ fås genom att maximera likelihood-funktionen

$$L(\alpha, \beta, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_1 - \alpha - \beta x_1)^2}{2\sigma^2}} \cdot \dots \cdot \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_n - \alpha - \beta x_n)^2}{2\sigma^2}} = (2\pi)^{n/2} \cdot (\sigma^2)^{n/2} \cdot e^{-\frac{1}{2\sigma^2} \sum (y_i - \alpha - \beta x_i)^2}$$

Hur än σ väljs så kommer L att maximeras med avseende på α och β då $\sum (y_i - \alpha - \beta x_i)^2$ är minimal, och eftersom det är just denna kvadratsumma som minimeras med MK-metoden så blir skattningarna av α och β de samma vid de två metoderna. Med ML-metoden kan vi dessutom skatta σ^2 varför vi väljer den. Logaritmeras likelihoodfunktionen fås

$$\ln L(\alpha, \beta, \sigma^2) = \frac{n}{2} \ln(2\pi) + \frac{n}{2} \ln(\sigma^2) - \frac{1}{2\sigma^2} \sum (y_i - \alpha - \beta x_i)^2$$

Deriveras denna med avseende på var och en av parametrarna och sedan sättes till noll fås ekvationssystemet

$$\frac{\partial \ln L}{\partial \alpha} = \frac{1}{\sigma^2} \sum (y_i - \alpha - \beta x_i) = 0 \quad (4)$$

$$\frac{\partial \ln L}{\partial \beta} = \frac{1}{\sigma^2} \sum (y_i - \alpha - \beta x_i)x_i = 0 \quad (5)$$

$$\frac{\partial \ln L}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum (y_i - \alpha - \beta x_i)^2 = 0 \quad (6)$$

att lösa med avseende på α , β och σ^2 . Eftersom vi kan förlänga de två första ekvationerna med σ^2 och därmed bli av med den kan vi använda dessa till att skatta α och β . (4) och (5) kan formas om till

$$\begin{aligned}\sum y_i &= n\alpha + \beta \sum x_i \\ \sum x_i y_i &= \alpha \sum x_i + \beta \sum x_i^2\end{aligned}\quad (7)$$

Delas första ekvationen med n fås

$$\bar{y} = \alpha + \beta\bar{x} \iff \alpha = \bar{y} - \beta\bar{x} \quad (8)$$

som vi kan stoppa in i (7) som då blir

$$\begin{aligned}\sum x_i y_i &= \bar{y} \sum x_i - \beta\bar{x} \sum x_i + \beta \sum x_i^2 \iff \\ \sum x_i y_i &= \beta(\sum x_i^2 - \bar{x} \sum x_i) + \bar{y} \sum x_i \iff \\ \beta &= \frac{\sum x_i y_i - \bar{y} \sum x_i}{\sum x_i^2 - \bar{x} \sum x_i} = \frac{\sum x_i (y_i - \bar{y})}{\sum x_i (x_i - \bar{x})} = [(2)] = \frac{\sum (x_i - \bar{x}) y_i}{\sum x_i (x_i - \bar{x})} = [(2) \text{ och } (3)] = \frac{S_{xy}}{S_{xx}}\end{aligned}\quad (9)$$

Detta resultat tillsammans med (8) ger ML-skattningarna av α och β

$$\beta^* = \frac{S_{xy}}{S_{xx}}, \quad \alpha^* = \bar{y} - \beta^* \bar{x}$$

Dessa värden insatta i (6) förlängd med σ^4 ger

$$(\sigma^2)^* = \frac{1}{n} \sum (y_i - \alpha^* - \beta^* x_i)^2$$

som dock inte är väntevärdesriktig utan korrigeras till

$$(\sigma^2)^* = s^2 = \frac{1}{n-2} \sum (y_i - \alpha^* - \beta^* x_i)^2 = \frac{Q_0}{n-2}$$

som är det. Q_0 som är summan av kvadratiska avvikelser från observationerna y_i till motsvarande punkt på den skattade linjen kallas *residualkvadratsumma* och den kan skrivas på formen

$$Q_0 = S_{yy} - \frac{S_{xy}^2}{S_{xx}}$$

7.3 Skattningarnas fördelning

Om vi börjar med β^* och utgår från (9)

$$\beta^* = \frac{S_{xy}}{S_{xx}} = \frac{\sum (x_i - \bar{x}) y_i}{\sum x_j (x_j - \bar{x})} = \sum c_i y_i \quad \text{där} \quad c_i = \frac{x_i - \bar{x}}{S_{xx}} \quad (10)$$

den är alltså en linjär funktion av de normalfördelade observationerna och därmed är skattningen normalfördelad. Väntevärdet blir

$$\begin{aligned}E(\beta^*) &= E\left(\sum c_i Y_i\right) = \sum c_i E(Y_i) = \sum c_i (\alpha + \beta x_i) = \frac{1}{S_{xx}} \sum (x_i - \bar{x})(\alpha + \beta x_i) \\ &= \frac{\alpha}{S_{xx}} \sum (x_i - \bar{x}) + \frac{\beta}{S_{xx}} \sum (x_i - \bar{x}) x_i = 0 + \beta \frac{S_{xx}}{S_{xx}} = \beta\end{aligned}$$

där vi i näst sista ledet åter använde hjälpresultaten (2) och (3). Skattningen är alltså väntevärdesriktig och dess varians blir

$$V(\beta^*) = V\left(\sum c_i Y_i\right) = \sum c_i^2 V(Y_i) = \sum c_i^2 \sigma^2 = \frac{\sigma^2}{S_{xx}^2} \sum (x_i - \bar{x})^2 = \frac{\sigma^2}{S_{xx}}$$

dvs

$$\beta^* = \frac{S_{xy}}{S_{xx}} \text{ är en observation av } \beta^* \in N\left(\beta, \frac{\sigma}{\sqrt{S_{xx}}}\right)$$

$\alpha^* = \bar{y} - \beta^*\bar{x}$ är även den normalfördelad eftersom den är en linjär funktion av normalfördelningar. Väntevärdet blir

$$\begin{aligned} E(\alpha^*) &= E(\bar{Y}) - \bar{x}E(\beta^*) = E\left(\frac{1}{n} \sum Y_i\right) - \bar{x}\beta = \frac{1}{n} \sum (\alpha + \beta x_i) - \bar{x}\beta = \\ &= \frac{1}{n} \sum \alpha + \frac{\beta}{n} \sum x_i - \bar{x}\beta = \alpha + \beta\bar{x} - \bar{x}\beta = \alpha \end{aligned}$$

så även α^* är väntevärdesriktig. Innan vi beräknar dess varians har vi nytta av att \bar{Y} och β^* är oberoende av varandra. Vi visar här att de är okorrelerade, vilket räcker för variansberäkningen. Återigen visar det sig fördelaktigt att uttrycka β^* enligt (10)

$$\begin{aligned} C(\bar{Y}, \beta^*) &= C\left(\frac{1}{n} \sum Y_i, \sum c_j Y_j\right) = \frac{1}{n} \sum_i \sum_j c_j C(Y_i, Y_j) = [Y_i \text{ är ober. av } Y_j \text{ då } i \neq j] = \\ &= \frac{1}{n} \sum c_i C(Y_i, Y_i) = \frac{1}{n} \sum c_i V(Y_i) = \frac{\sigma^2}{n} \sum c_i = \frac{\sigma^2}{nS_{xx}} \sum (x_i - \bar{x}) = 0 \end{aligned}$$

där vi återigen känner igen (1) i sista steget. Variansen för α^* blir

$$V(\alpha^*) = V(\bar{Y} - \beta^*\bar{x}) = V(\bar{Y}) + \bar{x}^2 V(\beta^*) - 2\bar{x}C(\bar{Y}, \beta^*) = \frac{\sigma^2}{n} + \bar{x}^2 \frac{\sigma^2}{S_{xx}} + 0$$

dvs

$$\alpha^* = \bar{y} - \beta^*\bar{x} \text{ är en observation av } \alpha^* \in N\left(\alpha, \sigma\sqrt{\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}}\right)$$

α^* och β^* är dock inte oberoende av varandra. Kovariansen mellan dem är

$$C(\alpha^*, \beta^*) = C(\bar{Y} - \beta^*\bar{x}, \beta^*) = C(\bar{Y}, \beta^*) - \bar{x}C(\beta^*, \beta^*) = 0 - \bar{x}V(\beta^*) = -\bar{x}\frac{\sigma^2}{S_{xx}}.$$

För variansskattningen och residualkvadratsumman gäller

$$(\sigma^2)^* = s = \frac{1}{n-2} \sum (y_i - \alpha^* - \beta^*x_i)^2 = \frac{Q_0}{f}, \quad \frac{Q_0}{\sigma^2} \in \chi^2(f)$$

Sakregister

- AR(1)-process, 32–35
 - minnesfaktor, 32
 - simulering av, 34
- autokorrelation, 28
- autokorrelationsfunktion, 28
 - Matlabkommandon, 30
 - skattning av, 29
- beroende variabel, 3
- fördelningsfritt test, 37
- förklarande variabel, 3
- förklaringsgrad, 15–16
- glidande medelvärde, 27
- jämförelse av två lutningar, 17–19
- kalibreringsintervall, 12
- Konfidensintervall, förväntat värde, 10
- konfidensintervall, förväntat värde, 9
- korrelation
 - anknytning till förklaringsgrad, 25
 - anknytning till linjär regression, 25
 - test av samband, 23–24
- korrelationsanalys, 22–25
- korrelationskoefficient
 - ”fallgropar”, 24
 - skattning, 23
 - tolkning, 24
- Mann-Kendalls test, 37–41
 - seasonal Kendall test, 39–41
 - skattning av trend, 38
- multipl linjär regression, 20–21
 - kovariansmatris, 20
 - matrisnotation, 20
 - parameterskattning med MK, 20
- oberoende variabel, 3
- orsakssamband, 24
- outliers, 16–17
- prediktionsintervall, observationer, 10–11
- prognosintervall, observationer, 10–11
- regression
 - enkel linjär, 5–19
 - konfidensintervall för parametrarna, 8–9
 - modellantaganden, 6–7
 - multipl linjär, 20–21
 - parameterskattningar, 7–8
 - parameterskattningarnas fördelning, 44
- regressionsfunktion, 3
- regressionslinje, 6
- residualanalys, 12–15
- residualer, 13
- responsvariabel, 3
- statistiskt samband, 24
- test av oberoende, 30
- test av samband, 9, 23
- tidsserier, 26–35
 - autokorrelation, 28
 - autokorrelationsfunktion, 28
 - autokorrelationsfunktion, skattning av, 29
 - karaktäristiska egenskaper, 26
 - komponentuppdelning, 27
 - outliers, 37
 - saknade värden, 37
 - stationär, 27
 - syfte med analys, 26
 - värden under detektionsnivå, 37
- trendanalys
 - beroende data, 35
 - exempel på, 9
 - fördelningsfritt, 37
 - skattning av trend, 38
- variationsuppdelning
 - oförklarad variation, 15
 - total variation, 15
 - variation förklarad av modell, 15
- vitt brus, 32

HT 2012

Matematisk statistik
Matematikcentrum
Lunds universitet
Box 118, 221 00 Lund

<http://www.maths.lth.se/>