



Stockholms universitet

OBS! Läs noga igenom anvisningarna i tentamen, t.ex. hur du ska skriva svaren. Det är ditt ansvar som student att följa de anvisningar som ges.

NOTE! Read the examination instructions carefully, e.g. how to write the answers. It is your responsibility as a student to follow the given instructions.

Skriv din anonymiseringskod och dagens datum på allt material du lämnar in.
(Enter your anonymization code and today's date on all submitted materials)

Anonymiseringskod (Anonymization code)	3	1	1	-	0	0	4	0	-	Z	H	G
Datum (Date YYYY-MM-DD)	2021-11-29							Plats nr. (Seat No.)	12			

Kurs/Kurskod (Course/Course code)	ST123G
Kursmoment (Course component)	Regressions och tidsserieanalys

Fylls i av tentamensvärd (To be filled in by invigilator)

Direkt i skrivning: (kryss)		Svarsblankett: (kryss)		Lösa svarsblad: (antal)	11.
--------------------------------	--	---------------------------	--	----------------------------	-----

Lämnat in blankt: (kryss)		Dator: (kryss)	
------------------------------	--	-------------------	--

Inlämningstid: 17:52 Signatur tentamensvärd: H/S

Fylls i av lärare/examinator (To be filled in by teacher/examinator)

Betyg:	A	Poäng:	90
--------	---	--------	----

29 28 7 15 11

Signatur rättande lärare/examinator: M



1. a) Vi behöver följande summor:

$$\sum x_i = 8$$

$$\sum y_i = 52$$

$$\sum x_i y_i = 60$$

$$\sum x_i^2 = 12$$

$$\sum y_i^2 = 358$$

$$n = 8$$

Med minsta kvadrat metoden kan vi få fram

$$b_1 = \frac{n \sum x_i y_i - (\sum x_i)(\sum y_i)}{n \sum x_i^2 - (\sum x_i)^2} = \frac{64}{32} = 2 \quad \text{(Skattning av } \beta_1 \text{)}$$

Med b_1 kan vi få fram a

$$a = \bar{y} - b_1 \bar{x} = 6,5 - 2 \cdot 1 = 4,5 \quad \text{(Skattning av } \beta_0 \text{)}$$

b) b_1 's koefficient säger att en ökning med en timmas fysisk aktivitet ökar oavbruten sömn med två timmar.

Det måste dock tas med försiktighet, vi har visat en korrelation, men inte kausalitet.

Det kan vara så att de som sover längre och tränar mer

a gör oss värde sömn en individ som mängden

inte utför någon fysisk aktivitet uppskattas ha.

I många fall ger interceptet ingen särskilt användbar information, om datat inte innehåller x-värden med noll t.ex. (extrapolering)

I detta fall gör datat det, så man kan säga att de som inte tränar sover 4,5 timmar.

Uppg.nr.: (Task no.)

Lärarens kommentar: (Teacher's note)

Poäng: (Points)

Uppg.nr.:
(Task no.)

Lärarens
kommentar:
(Teacher's
note)

Poäng:
(Points)



1. b) fort. Vanligtvis vill man också testa koefficienterna innan man hävdar att det råder något samband.

Uppg.nr.:
(Task no.)

Lärens kommentar:
(Teacher's note)

1. c) σ_E^2 är variansen för feltermerna. Den är ett mått på det teoretiska "bruset" som gör att inte ens vår teoretiska modell kan ge perfekta prediktioner. Våra modeller ställer särskilda krav på σ_E^2 's egen skapar, t.ex. att den är normalfördelad med väntevärde noll. Att feltermerna är oberoende av varandra, och att variansen är konstant. (Homoskedastisk)

Tyvärr är dessa feltermar omöjliga att komma åt, så vårt enda sätt att undersöka dem är att använda våra residuer för att bilda en skattning.

Vi skattar σ_E^2 med S_E^2 .

$$S_E^2 = \frac{\sum (y_i - \hat{y}_i)^2}{n-k-1} = \frac{4}{6} = \frac{2}{3}$$

Poäng:
(Points)

Uppg.nr.:
(Task no.)

Lärarens
kommentar:
(Teacher's
note)

Poäng:
(Points)



1.d)

$$s_{b_1}^2 = \frac{s_e^2}{\sum (x_i - \bar{x})^2} = \frac{2/3}{4} = \frac{1}{6}$$

$$s_{b_1} = \sqrt{\frac{1}{6}} \approx 0,40825$$

Det visar hur mycket b_1 kan variera.
 Om s_{b_1} är väldigt hög gör det skattningen mer osäker och konfidensintervallet blir större.

1.e) Jag väljer att testa b_1 . (T-test)

Hypoteser $\alpha = 0,05$, tvås. d. st test så använder
 $(H_0: \beta_1 = 0)$
 $(H_A: \beta_1 \neq 0)$ $\frac{\alpha}{2} = 0,025$

$$\frac{b_1 - \beta_0}{s_{b_1}} = \frac{2}{\sqrt{1/6}} \approx 4,89898$$

$$T_{obs} = 4,89898 > T_{krit} = T_{8-1, \frac{0,05}{2}} = T_{6, 0,975} = 2,447$$

Förkasta H_0 , b_1 är signifikant under en 5% signifikansnivå.

Det betyder att ~~ett~~ timmar av fysisk aktivitet har en effekt på timmar sömn, med en 5% risk att sambandet ej gäller och bara visades av ren slump.

Vi kan inte härda kausaltet.

Uppg.nr.:
(Task no.)

Lärares
kommentar:
(Teacher's
note)

Poäng:
(Points)



2.a)

T-test $\frac{\alpha}{2} = 0,025$

Uppg.nr.: (Task no.)

Hypoteser

$H_0: \beta_3 = 0$
 $H_A: \beta_3 \neq 0$
 $T_{obs} = \frac{0,64817 - 0}{0,48619} \approx 1,3332$

Krit (approx. till 2)

$T_{obs} \approx 1,3332 < T_{krit} = T_{n-k-1} = T_{100-3-1} = T_{96, 0,975} \approx 1,96$

Vi kan ej förkasta H_0 , med vår signifikansnivå är β_3 insignifikant.

Lärarens kommentar: (Teacher's note)

Tolkning av koefficienten:

Balkong är en dummyvariabel, när den är lika med ett ökar \hat{y} med 0,64817. Man kan se det som interceptet ökar med summan av koefficienten.

2.b) KI: $b_1 \pm T_{96, 0,975} \cdot s_{b_1}$

Använder T för mer precisa siffror

Ⓐ $0,05978 - 1,985 \cdot s_{b_1} = 0,03994$

$s_{b_{1A}} = 0,0099949622$

Ⓑ $0,05978 + 1,985 \cdot s_{b_1} = 0,07961$

$s_{b_{1B}} = 0,0099899244$

approx. s_{b_1}

$\frac{s_{b_{1A}} + s_{b_{1B}}}{2} = 0,0099924433$

Hypoteser

T-test, $\frac{\alpha}{2} = 0,025$

$H_0: \beta_1 = 0$

$H_A: \beta_1 \neq 0$

$T_{obs} = \frac{0,05978 - 0}{0,0099924433} \approx 5,9825$

$T_{obs} \approx 5,9825 > 1,985 = T_{krit}$

Förkasta H_0 , b_1 signifikant.

Summan statsats kan dras direkt från "95% confi. lim". Eller som den utnyttjar en 5% signifikans för att skapa \rightarrow

Poäng: (Points)

Uppg.nr.:
(Task no.)

Lärarens
kommentar:
(Teacher's
note)

Poäng:
(Points)



2.b) föresättning.

→ intervallet. Eftersom intervallet inte innehåller 0 vet vi med 95% konfidens att b_1 är mellan 0,03994 och 0,07961.

2c) vi ser att rum och kvm har ett högt VIF, över 10.

Dom har alltså ett starkt samband, vilket är intuitivt. Ett problem som kan uppstå i sådana fall är att koefficienter blir överskattade. Här ser det ut som att antalrum är negativ pga att b_1 är överskattad, och b_2 "balanserar".

Man ska också komma ihåg att b_2 är insignifikant, så koefficienten bör förkastas.

Regressionen förbättras nog av att ta bort antingen kvm eller antalrum, eftersom dom är korrelerade

kommer en beskriva en stor del av bådars påverkan.

Förmodligen är det ändå bättre att välja bort antalrum.

Uppg.nr.: (Task no.)

Lärarens kommentar: (Teacher's note)

Poäng: (Points)

Uppg.nr.:
(Task no.)

Lärarens
kommentar:
(Teacher's
note)

Poäng:
(Points)



2.d) $R^2 = \frac{SSR}{SST} = \frac{SST - SSE}{SST} = \frac{582,2127}{732,29457} = 0,79505$

Detta modellens förklaringsvärde, men ~~en~~ R^2 kan ökas av att lägga till onödiga koefficienter.

Därför används ofta R^2_{adj} för multipel regression.

$R^2_{adj} = 1 - \frac{MSE}{SST/(n-1)} = 1 - \frac{1,56335}{7,396915} \approx 0,211352$

Alltså beskriver modellen ^{ca} 21,1352% av variationen av y. Vi ser här en väldigt skillnad i R^2 och R^2_{adj} .

2.e) F-test med 5% sign. $\alpha = 0,05$

Hypoteser

$H_0: \beta_1 = \beta_2 = \beta_3 = 0$
 $H_A: \beta_1 \text{ eller } \beta_2 \text{ eller } \beta_3 \neq 0$

$F_{obs} = \frac{MSR}{MSE} \approx 124,138$

$F_{obs} = 124,138 \rightarrow F_{crit} = F_{k, n-k-1} = F_{3, 96} \approx 2,7$
 $\alpha = 0,05$

Förkasta H_0 , åtminstone en av $\beta_1, \beta_2, \beta_3$ är signifikant på en signifikansnivå på 5%.

Uppg.nr.: (Task no.)

Lärarens kommentar: (Teacher's note)

Poäng: (Points)

Uppg.nr.:
(Task no.)

Lärarens
kommentar:
(Teacher's
note)

Poäng:
(Points)

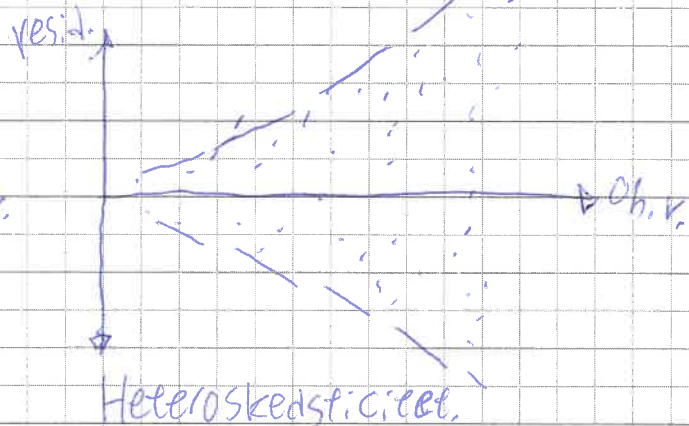
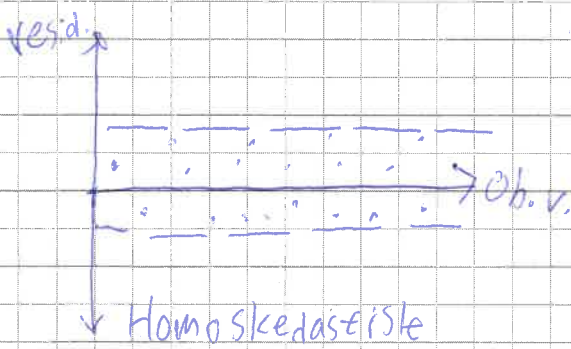


3.a) Variansen är konstant när den beroende variabeln y förändras

Uppg.nr.: (Task no.)

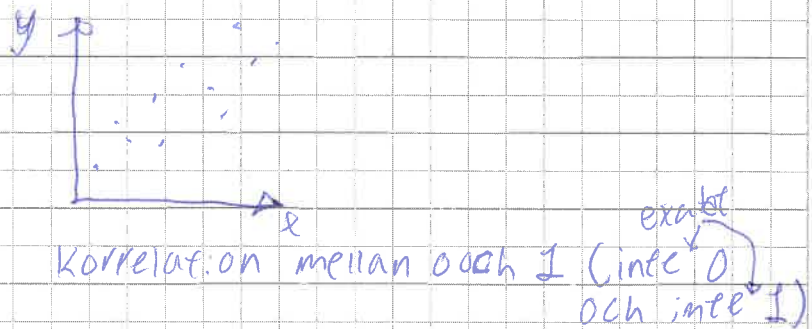
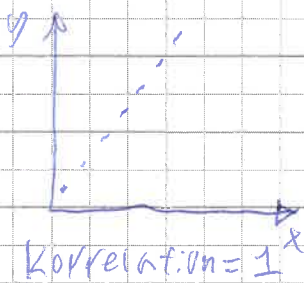
Oberoende variabel och residualer

Lärarens kommentar: (Teacher's note)



3.b) Mäter linjärt samband, men inte lutning.

En beräkning på styrkan av det linjära sambandet.



3.c) Test för tidsserier

Autokorrelation?

Ska vara mellan $\frac{2}{n}$ och $\frac{1}{n}$

Uppg.nr.:
(Task no.)

Lärarens
kommentar:
(Teacher's
note)

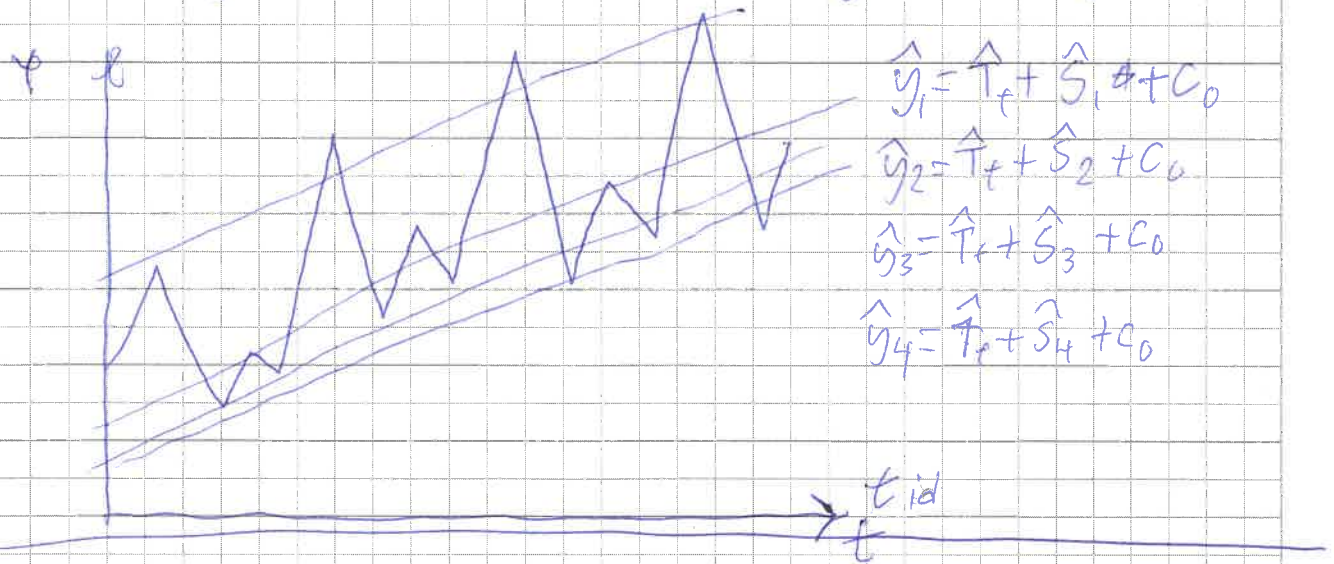
Poäng:
(Points)



3.d) En uppskattning av trenden i tidsserier.

~~Även värtier per säsong~~

Ändra den trenden för att räkna ut $y_t - \hat{T}_t$.
Säsongss indelar $y_t - \hat{T}_t$ och tar medelvärde för
säsongen för att få en skattning av säsongseffekten.



Uppg.nr.:
(Task no.)

Lärarens
kommentar:
(Teacher's
note)

3.e) Till för att välja antal ob. variabler.

Lägger en restriktion på antal ob. v. λ

Bra för ~~data set~~ enorma data set med massa
variabler. Också bra för att välja

polynomgrad, ALLTÄR $x_1^2, x_1^3, \dots, x_1^{27}$ t.ex..

Poäng:
(Points)

Uppg.nr.:
(Task no.)

Lärarens
kommentar:
(Teacher's
note)

Poäng:
(Points)



4.a) En autokorrelationsfunktion visar hur mycket

y_t korrelerar med y_{t-1} ; ~~ja~~ man kan kolla

korrelationen så länge tillbaka som datat tillåter.

~~$\gamma(k)$~~

~~Man kan ha fler jägande y värden
samt y_{t-1}~~

~~Från figuren ser vi att konf. intervall skapas~~

~~med vår första lagning, y_{t-1} . Och att det ökar~~

~~efter den fjärde, y_{t-4} . Jag tolkar detta som att~~

~~en ökning i längden för KI innebär ett mindre
höjden~~

~~KI, alltså osäkerhetsprognos~~

y mot y har en korrelation på 1, den minskar

sedan men ökar vid y mot y_{t-4} , det är ett

tecken på att det finns en säsongs effekt ~~och~~ var

fjärde kvarten.

Osäkerheten i korrelationen ser ut att öka längre

bak i tiden, pga KI.

Uppg.nr.:
(Task no.)

Lärens
kommentar:
(Teacher's
note)

Poäng:
(Points)

Uppg.nr.:
(Task no.)

Lärarens
kommentar:
(Teacher's
note)

Poäng:
(Points)



4.b) Eftersom det ser ut att finnas en säsongseffekt är det nog bäst att ta modell AR(4).

Här kan vi också se att y_{2021K1} har en signifikant koefficient. Vi ser också i figuren att v_t inte

har något tydligt samband hos residualerna. Det är

ett tecken på att vi har lyckats säsongseffekt
anpassa modellen

$$4.c) \hat{y}_{2021K1} = 62,2752 + 0,3321 \cdot y_{2020K4} + 0,0137 \cdot y_{2020K3}$$

$$\hat{y}_{2021K1} = 96,852102$$

$$\hat{y}_{2021K2} = 62,2752 + 0,3321 \cdot \hat{y}_{2021K1} + 0,0137 \cdot y_{2020K4}$$

$$\hat{y}_{2021K2} = 95,80772807$$

Uppg.nr.:
(Task no.)

Lärarens
kommentar:
(Teacher's
note)

Poäng:
(Points)



Uppg.nr.:
(Task no.)

Lärarens
kommentar:
(Teacher's
note)

5.a)
$$OR(x_i) = \frac{0,0087395201}{0,0084338661} = 1,036241267$$

Ålder ökar sannolikheten ~~att~~ få diabetes med ungefär 3,6% per år.

5.b) Logaritmen av en logistisk regression blir linjär. Är icke linjär då man använder den i sin exponentiella form. Beräknas ofta med maximum likelihood

5.c)
$$e^{(-5,0009 - 0,0101 \cdot 100 + 0,1062 \cdot 40 + 0,0356 \cdot 60 + 0,0937 \cdot 4)} = 2,112558981$$

$$P(\text{diabetes}) = \frac{2,112558981}{3,112558981} = 0,6787209477$$

Poäng:
(Points)

Uppg.nr.:
(Task no.)

Lärarens
kommentar:
(Teacher's
note)

Poäng:
(Points)