



EKSAMENSOPPGAVE

Institutt: IKBM

Eksamen i: STAT 100 STATISTIKK

Tid: Torsdag 11. desember 9:00 – 12:30 (3.5 timer)

Emneansvarlig: Solve Sæbø, Tlf 67232561

Tillatte hjelpemidler:

C3: alle typer kalkulatorer, alle andre hjelpemidler

Oppgaveteksten er på: 12 (inkludert svarskjema)
antall sider inkl. vedlegg

Alle deloppgaver teller likt. For hver oppgave er det 5 svaralternativer. Kun ett svaralternativ er riktig. Du får ett poeng for riktig svar, null poeng for feil svar. Maksimal score er da 40 poeng.

Alle svar føres i svarskjemaet på side 12, og denne siden er den ENESTE som skal leveres når eksamen er slutt. Husk å skrive kandidatnummer på svarskjemaet!

Sammenlikning av bakteriegener (Oppgave 1-7)

Tre av de vanligste bakterieslektene i jord er *Pseudomonas*, *Dechloromonas* og *Paracoccus*. I en studie plukket man ut et lite fragment av et gen som alle tre slektene har (en kort DNA-bit), og man fant sekvensen av disse (rekkefølgen av basene: A, T, C og G). For eksempel fant man at genbiten for *Pseudomonas* hadde følgende sekvens med 41 baser:

GAAGGGGTGGCACCGCCCTGGTGTGGCGCTGAAATATACCA

En måte å sammenlikne bakterier på er å telle opp hvor mange bokstaver man har av hver type i et gen. Tabellen nedenfor angir en opptelling av basene i DNA-bitene for de tre slektene:



	A	T	C	G	sum
Pseudomonas	9	7	10	15	41
Dechloromonas	4	2	14	15	35
Paracoccus	1	5	15	17	38
sum	14	14	39	47	114

Man ønsket å teste om basefordelingen var avhengig av bakterieslekt og utførte en kji-kvadrattest i R Commander med følgende utskrift (noen tall er erstattet med «?»)

```
> round(.Test$expected,2) # Expected Counts
      A      T      C      G
Pseudomonas  ? 5.04 14.03 16.90
Dechloromonas 4.30 4.30 11.97 14.43
Paracoccus    4.67 4.67 13.00 15.67

> round(.Test$residuals^2, 2) # Chi-square Components
      A      T      C      G
Pseudomonas  3.12 0.77 1.16 0.21
Dechloromonas 0.02 1.23 0.34 0.02
Paracoccus    ? 0.02 0.31 0.11

> print(.Test)
      Pearson's Chi-squared test

data:  .Table
X-squared =  ? , df = ?, p-value = 0.1165
```

Bruk resultatene i den grad du finner det nødvendig for å svare på oppgavene nedenfor.

Oppgave 1

Hvor mange frihetsgrader er tilknyttet den kji-kvadratfordelte testobservatoren Q i testen som er beskrevet ovenfor?

- a) 2 b) 3 c) 12 d) 6 e) 9

Oppgave 2

Dersom en nullhypotese om at basefordelingen er uavhengig av bakterieslekt er sann, hva er da forventet antall av basen «A» i sekvensen til Pseudomonas?

- a) 4.92 b) 4.32 c) 4.67 d) 9.12 e) 5.04

Oppgave 3

Hva er bidraget (Chi-square component) til testobservatoren Q fra kombinasjonen «A» og slekten Paracoccus?

- a) 2.88 b) 4.12 c) 1.18 d) 0.79 e) 3.01

Oppgave 4

Dersom man har lite data, kan man lett gjøre en såkalt type II feil i en slik hypotesetestingssituasjon. Hva er en type II feil?

- a) Det betyr at forventet antall under nullhypotesen er estimert feil.
b) Det betyr at man beholder nullhypotesen selv om man burde ha forkastet den.
c) Det betyr at den alternative hypotesen er feil.
d) Det betyr at man forkaster nullhypotesen selv om den er sann.
e) Det betyr at hypotesene er satt opp på feil premisser.



Oppgave 5

Og hva er en type I feil?

- Det betyr at forventet antall under nullhypotesen er estimert feil.
- Det betyr at man beholder nullhypotesen selv om man burde ha forkastet den.
- Det betyr at null-hypotesen er feil.
- Det betyr at man forkaster nullhypotesen selv om den er sann.
- Det betyr at hypotesene er satt opp på feil premisser.

Av biologiske grunner så er det mer korrekt å telle «A» og «T» sammen som én kategori i en oppdeling av basene. Det samme gjelder «C» og «G». Derfor slo man disse kolonnene sammen i to nye kategorier «A/T» og «C/G» og kjørte en ny analyse med følgende utskrift:

```
> .Table # Counts
      A/T C/G
Pseudomonas  16 25
Dechloromonas  6 29
Paracoccus    6 32

> round(.Test$expected,2) # Expected Counts
      A/T  C/G
Pseudomonas 10.07 30.93
Dechloromonas 8.60 26.40
Paracoccus   9.33 28.67

> round(.Test$residuals^2, 2) # Chi-square Components
      A/T  C/G
Pseudomonas 3.49 1.14
Dechloromonas 0.78 0.26
Paracoccus  1.19 0.39

> print(.Test)
      Pearson's Chi-squared test

data:  .Table
X-squared = 7.2463, df = 2, p-value = ?
```

Oppgave 6

Hva er det laveste av følgende testnivå som gir forkastning av en nullhypotese om at basefordelingen er uavhengig av bakterieslekt i denne testen?

- 0.1
- 0.05
- 0.025
- 0.01
- 0.005

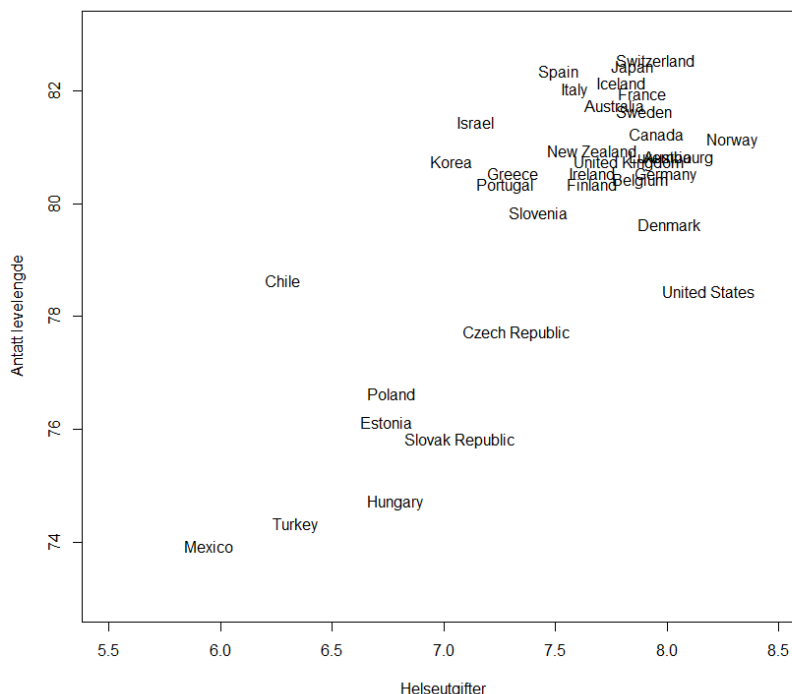
Oppgave 7

Av hvilken statistisk grunn er den siste testen med to base-kategorier mer pålitelig enn den første med fire basekategorier?

- Fordi verdien på Q er nærmere 0.
- Fordi vi har færre antall frihetsgrader.
- Fordi det ble større utslag på enkeltbidragene til testobservatoren Q (Chi-square components).
- Fordi antagelsen om kji-kvadratfordelt Q er mer korrekt i den siste testen på grunn av høyere forventet antall i hver celle i tabellen.
- Fordi p-verdien ble lavere (p-verdien er ikke gjengitt i siste test, men den ble lavere).

Levelengde og helseutgifter (Oppgave 8-16)

Fra OECD sin database har vi hentet ut antatt levelengde for barn født i 2011 og et mål på totale offentlige helseutgifter pr innbygger i 33 OECD medlemsland for samme år. Disse variablene er plottet mot hverandre i Figur 1.



Figur 1: Antatt levealder for nyfødte i 2011 i ulike land plottet mot et samlet mål på totale offentlige helseutgifter i de samme landene. Kilde: OECD.StatExtracts.

En regresjonsanalyse ble kjørt i R Commander med levelengde som respons Y , og helseutgifter som forklaringsvariabel X . Modellen som antas er $Y_i = \alpha + \beta X_i + \epsilon_i$ der $\epsilon_i \sim N(0, \sigma)$. Resultatet av analysen ble:

```
> summary(LinearModel.1)

Call:
lm(formula = Y ~ X, data = data2011)

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  56.5220     3.7399  15.113 7.55e-16 ***
X              3.1463     0.4984    ?      ?

s: 1.654 on 31 degrees of freedom
Multiple R-squared: ? ,
F-statistic: 39.85 on 1 and 31 DF, p-value: ?
```

(Fortsetter på neste side)



ANOVA-tabell:

	Df	Sum Sq	Mean Sq
X	1	109.04	109.04
Residuals	31	84.82	2.74

Oppgave 8

Vi skal se på testingen av hypotesene: $H_0 : \beta = 0$ mot $H_1 : \beta \neq 0$. Hvilket av følgende utsagn om *både* verdi av testobservator T for testen og konklusjon om sammenheng mellom helseutgifter og levealder stemmer overens med resultatene fra R Commander?

- $T = 6.31$, forkaster H_0 med testnivå 1% og hevder at det er sammenheng.
- $T = 3.15$, forkaster ikke H_0 med testnivå 1% og hevder at det ikke er sammenheng.
- $T = 6.31$, forkaster ikke H_0 med testnivå 1% og hevder at det ikke er sammenheng.
- $T = 3.15$, forkaster H_0 med testnivå 1% og hevder at det er sammenheng.
- $T = 6.31$, forkaster H_0 med testnivå 1% og hevder at det ikke er sammenheng.

Oppgave 9

Hva blir et 95% konfidensintervall for β ?

- [2.29, 3.99]
- [2.12, 4.16]
- [2.17, 4.11]
- [1.86, 4.42]
- [1.93, 3.97]

Oppgave 10

Hvordan tolkes et 95% konfidensintervall for β ?

- Det er 95% sannsynlig at $\hat{\beta}$ ligger i intervallet.
- Ved gjentatte trekninger fra samme populasjon vil 95% av beregnede intervaller inneholde den sanne verdien av β .
- Ved gjentatte trekninger fra samme populasjon vil 95% av beregnede intervaller inneholde den estimerte effekten av helseutgifter på levealder.
- Det er 95% sannsynlig at forventet levealder ligger i intervallet.
- Det er 95% sannsynlig at den estimerte linja ligger i intervallet.

Oppgave 11

Hva er estimert forventet levealder i OECD-land dersom det offentlige ikke bruker penger på helse? (Dette er nok et urealistisk situasjon, men vi spør likevel)

- 33.3
- 40.3
- 0
- 3.15
- 56.5

Oppgave 12

I henhold til den estimerte modellen vil et land som bruker $X=7$ på helse ha en predikert levealder på 78.5 år. Gjennomsnittlig helseutgifter i de 33 undersøkte landene i 2011 var $\bar{X} = 7.48$. Hva blir et 95% *prediksjonsintervall* for Y når $X=7$?

- [75.0, 82.0]
- [76.0, 81.0]
- [74.0, 83.0]
- [74.1, 82.9]
- [77.7, 79.3]



Oppgave 13

Hvordan tolker vi prediksjonsintervallet vi fant i forrige oppgave?

- I et vilkårlig land med gjennomsnittlige helseutgifter vil man med 95% sannsynlighet ha en levealder i intervallet.
- Det er 95% sannsynlig at helseutgiftene til et vilkårlig land ligger i intervallet.
- I et vilkårlig land med helseutgifter på $X=7$ vil man med 95% sannsynlighet ha en levealder i intervallet.
- For et stort antall land med helseutgifter på $X=7$ vil med 95% sannsynlighet den gjennomsnittlige levealderen ligge i dette intervallet.
- Det er 95% sannsynlig at gjennomsnittlige helseutgifter i OECD-land ligger i intervallet.

Oppgave 14

I Korea har man en offentlige helseutgifter på $X=7$ og en levealder på 81 år. Hva er residualet til denne observasjonen i henhold til den estimerte modellen?

- 3.5
- 3.5
- 2.5
- 1.5
- 1.5

Oppgave 15

Hvordan kan vi tolke parameteren σ i denne modellen?

- Den beskriver variasjon i helseutgifter blant land med samme levelengde for sine innbyggere.
- Den beskriver forventet levelengde i et tenkt land med ingen offentlige helseutgifter.
- Den beskriver den totale variasjonen i levelengder uavhengig av helseutgiftsnivå.
- Den beskriver effekten av helseutgifter på levelengde i OECD-land.
- Den beskriver variasjon i levelengder blant land med identiske offentlige helseutgifter.

Oppgave 16

Hvor stor er R^2 (determinasjonskoeffisienten) ifølge den estimerte modellen?

- 0.86
- 0.45
- 0.78
- 0.56
- 0.51

Levelengder i ulike år (Oppgave 17-21)

Fra OECD sin database har vi også fått antatt levelengde for barn født i 2010 og i 2011 i de 33 OECD-landene. En oppsummering av dataene er gitt i tabellen nedenfor. «Diff 2011-2010» er de parvise differansene mellom levelengde i 2011 og 2010 for hvert land.

	mean	sd	n
2010	79.78	2.50	33
2011	80.06	2.46	33
Diff 2011-2010	0.28	0.19	33

Korrelasjon mellom 2010 og 2011: $r= 0.997$

Oppgave 17

Vi skal først betrakte observasjonene fra 2010 og 2011 som to uavhengige utvalg og vil kjøre en to-utvalgs test for å sammenlikne forventningen i 2010 (μ_{2010}) med forventningen i 2011 (μ_{2011}). I en slik test antar vi som regel et felles standardavvik σ for begge utvalg.

Hva blir verdien av estimatoren S_p for σ , dvs «pooled» standardavvik?

- 2.48
- 4.96
- 0.04
- 2.47
- 3.51



Oppgave 18

En to-utvalgstest for å teste $H_0 : \mu_{2010} = \mu_{2011}$ mot et tosidig alternativ gav en p-verdi på 0.65. Hvorfor er det meget stor grunn til å tro at en parvis test vil gi mye lavere p-verdi?

- Fordi vi får halvert antall frihetsgrader.
- Fordi antagelsen om felles standardavvik for begge år, som man gjør i to-utvalgstesten, er feil.
- Fordi levelengdene i de to årene er sterkt korrelert.
- Fordi parvis test alltid gir lavere p-verdi.
- Fordi de observerte gjennomsnittene er så like sett i forhold til standardavvikene.

Oppgave 19

Hva blir verdien av testobservatoren for en parvis test for å teste $H_0 : \mu_{2010} = \mu_{2011}$ mot et tosidig alternativ?

- 17
- 4.7
- 8.5
- 3.3
- 0.26

Oppgave 20

Man har egentlig grunn til å tro at levealderen er på veg opp i OECD-landene. Hva blir så verdien av testobservatoren for en parvis test for å teste $H_0 : \mu_{2010} = \mu_{2011}$ mot det ensidige alternativet $H_1 : \mu_{2010} < \mu_{2011}$?

- 8.5
- 17
- 4.7
- 3.3
- 10.2

Oppgave 21

Den tosidige testen fra oppgave 19 gav en p-verdi på $1.0 \cdot 10^{-9}$ og klar signifikans for at forventningene i de to årene er ulike. Hva er riktig å si om p-verdien for testen med et ensidig alternativ, som beskrevet i oppgave 20?

- Vi får samme p-verdi som for tosidig test.
- Vi får dobbelt så stor p-verdi som for tosidig test.
- P-verdien vil bli noe større, men vanskelig å si nøyaktig uten statistisk programvare.
- Vi får halvparten så stor p-verdi som for tosidig test.
- P-verdien blir lik testnivået.

Inntektsnivå hos kunder av ulike mobiloperatører (Oppgave 22-33)

Vi vil sammenlikne inntektsnivået hos kunder av tre ulike mobilfirmaer. Inntekt (Y) (målt i antall tusen norske kroner) ble observert hos noen tilfeldig valgte kunder i hvert firma, og en ANOVA-modell ble brukt for å analysere dataene.

Modellen som er antatt er: $Y_{ij} = \mu_i + \epsilon_{ij}$, der $\epsilon_{ij} \sim N(0, \sigma)$, $i = 1, 2, 3$ og $j = 1, \dots, n_i$.

En R Commander utskrift er gitt nedenfor:

```
> summary(AnovaModel.2)
              Df Sum Sq Mean Sq F value Pr(>F)
Firma          2 124538   62269      ?  0.0202 *
Residuals     18 229376   12743
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

>
              mean      sd      n
Firma1    342.5    130.58     8
Firma2    167.1    113.54     7
Firma3    308.3     80.85     6
```



Oppgave 22

Hvor mange frihetsgrader er knyttet til SSG, dvs kvadratsummen til gruppevariabelen «Firma»?

- a) 2 b) 18 c) 20 d) 21 e) 3

Oppgave 23

En av størrelsene vi beregner i denne variansanalysen er $\sum_{i=1}^3 n_i(\bar{Y}_i - \bar{Y})^2$. Hva kaller vi denne?

- a) SST b) SSG c) MSG d) MSE e) SSE

Oppgave 24

Hva blir verdien av den F-fordelte testobservatoren som brukes til å teste en nullhypotese om at forventet inntekt hos kunder er den samme i alle tre mobilfirmaer?

- a) 0.53 b) 3.16 c) 7.76 d) 0.20 e) 4.89

Oppgave 25

La F_{obs} være det tallet du fant i oppgave 24. Hva er et annet navn på sannsynligheten $P(F > F_{obs})$?

- a) Testnivået.
b) Type II feilen.
c) Type I feilen.
d) p-verdien for testen.
e) Determinasjonskoeffisienten.

Oppgave 26

Hvilken av følgende funksjoner kan gjennom en kontrastanalyse brukes til å teste om forventningen til inntekt for kunder i Firma 2 er lavere enn forventningen til inntekt til kunder i Firma 1 og 3. Følgende hypoteser skal brukes i testen: $H_0 : \theta = 0$ mot $H_1 : \theta > 0$

- a) $\theta = \mu_1 - \mu_2 + \mu_3$
b) $\theta = (\mu_1 + \mu_2)/2 - \mu_3$
c) $\theta = \mu_1 - 2\mu_2 + \mu_3$
d) $\theta = (\mu_1 - \mu_2 + \mu_3)/3$
e) $\theta = \mu_1 - \mu_3$

Oppgave 27

Man ville også teste om forventningene i Firma 1 og 3 var forskjellige fra hverandre og definerte følgende kontrast for dette:

$$\theta = \frac{1}{2}\mu_1 - \frac{1}{2}\mu_3$$

Hva blir et forventningsrett estimat for θ ?

- a) 23.5 b) 34.2 c) 17.1 d) 8.5 e) 19.3

Oppgave 28

Og hva blir standardfeilen til $\hat{\theta}$ funnet i forrige oppgave?

- a) 19.96 b) 30.48 c) 60.96 d) 27.15 e) 23.12



Oppgave 29

Vi betrakter fortsatt kontrasten definert i oppgave 27 og vil teste $H_0 : \theta = 0$ mot $H_1 : \theta \neq 0$

Hvilken fordeling har testobservatoren definert ved:

$$\frac{\hat{\theta} - 0}{SE(\hat{\theta})}$$

- a) Standard normalfordelt.
- b) T-fordelt med 2 frihetsgrader.
- c) F-fordelt med 3 og 18 frihetsgrader.
- d) T-fordelt med 18 frihetsgrader.
- e) Kji-kvadratfordelt med $(n-1)(k-1)$ frihetsgrader.

Oppgave 30

En av kundene hos Firma 1 hadde en inntekt på 350 (tusen kroner). Hva er residualet til denne observasjonen?

- a) 2.0
- b) 130.6
- c) 11.5
- d) 7.5
- e) 15.3

Oppgave 31

Hva er populasjonsstandardavviket σ et mål på?

- a) Variasjon i inntekt mellom kunder av samme mobilfirma.
- b) Variasjon i forventet inntekt mellom kunder av de tre firmaene.
- c) Forventet inntektsnivå hos kunder i samme firma.
- d) Generell variasjon i inntekt hos mobilkunder på tvers av mobilfirmaer.
- e) Variasjon mellom firmaer med hensyn til kunders inntektsnivå.

Oppgave 32

Noen hevder at forventet inntekt hos kunder av Firma 1 (μ_1) er lavere enn 400 (tusen) og ber deg teste dette med en hypotesetest. Hvilket par av hypoteser bør du da bruke?

- a) $H_0 : \mu_1 = 400$ mot $H_1 : \mu_1 < 400$
- b) $H_0 : \mu_1 = 0$ mot $H_1 : \mu_1 < 0$
- c) $H_0 : \mu_1 = 400$ mot $H_1 : \mu_1 > 400$
- d) $H_0 : \mu_1 = 400$ mot $H_1 : \mu_1 \neq 400$
- e) $H_0 : \bar{Y}_1 = 400$ mot $H_1 : \mu_1 < 400$

Oppgave 33

Hvis du vil utføre testen som er definert i oppgave 32, hvilken verdi får da den t-fordelte testobservatoren du bør bruke (hvis du på tror modellen som ble definert i innledningen før oppgave 22)?

- a) -4.12
- b) -2.12
- c) -1.44
- d) -0.78
- e) -1.67



Bilsalg (Oppgave 34-40)

Et bilfirma har $n=12$ salgsdistrikter med varierende antall selgere (utsalgssteder). Gjennomsnittlig salg pr uke (Y) og antall selgere (X) ble observert i alle distrikter over en periode på 6 måneder. Nedenfor er noen resultater fra en analyse i R Commander, blant annet en anovatabell fra en regresjonsanalyse.

	mean	sd	n
Salg (Y)	16.33	2.42	12
Selgere(X)	5.92	1.44	12

Korrelasjon mellom Y og X: $r = 0.55$

```
> Anovatabell
```

	Df	Sum Sq	Mean Sq
Selgere	1	19.859	19.859
Residuals	10	44.807	4.481

Anta en regresjonsmodell for dataene: $Y_i = \alpha + \beta X_i + \epsilon_i$, der $\epsilon_i \sim N(0, \sigma)$ for $i = 1, \dots, 12$

Oppgave 34

Hvordan ville du ha forklart parameteren β i modellen til noen som ikke kan statistikk?

- Den angir endringen i antall selgere som trengs for i gjennomsnitt å selge én bil mer pr uke.
- Den er endringen i forventningen til antall biler solgt (i gjennomsnitt) pr uke dersom antall selgere øker med én.
- Den er lik standardavviket i gjennomsnittsantall biler som selges, dersom antall selgere holdes konstant.
- Det er forventet gjennomsnittssalg når antall selgere er lik 0.
- Dersom antall selgere økes med én, vil β være lik økningen i antall solgte biler for en vilkårlig valgt bilforretning.

Oppgave 35

Hva er minste kvadraters estimatet for β i modellen?

- a) 0.55 b) 0.33 c) 5.92 d) 0.30 e) 0.92

Oppgave 36

Hva er et omtrentlig estimat på forventet salg pr uke dersom $X = 6$, altså et estimat på $E(Y|X = 6)$? (Legg merke til at ifølge R-utskriften er gjennomsnittet $\bar{X} \approx 6$)

- a) 10 b) 20 c) 6 d) 18 e) 16

Oppgave 37

Minste kvadraters estimatoren for β er en forventningsrett estimator. Hva menes med dette?

- Estimatoren gir alltid det man forventer å få.
- Estimatoren gir den rette forventningen til Y for gitt verdi av X .
- I det lange løp vil denne estimatoren i gjennomsnitt bli lik den sanne forventningen til X .
- I det lange løp vil denne estimatoren i gjennomsnitt bli lik β .
- I det lange løp vil denne estimatoren ha minst varians blant alle estimatorer for β .



Oppgave 38

Hva er et forventningsrett estimat for σ^2 , dvs populasjonsvariansen til støyleddet i modellen?

- a) 5.86 b) 4.48 c) 4.43 d) 2.07 e) 2.42

Oppgave 39

I regresjonsanalysen har vi STAT100 gjort en del modellantagelser som blant annet involverer responsvariabelen Y og forklaringsvariabelen X . Hvilken av følgende antagelser har vi ikke gjort?

- a) Vi har antatt at $X_i \sim N(\mu, \sigma)$.
b) Vi har antatt at σ er uavhengig av verdien til X .
c) Vi har antatt at $Y_i \sim N(\alpha + \beta X_i, \sigma)$.
d) Vi har antatt at alle støyleddene er uavhengige av hverandre.
e) Vi har antatt at støyleddene er normalfordelte.

Oppgave 40

I enkel regresjonsanalyse er antall frihetsgrader tilknyttet MSE lik $n - 2$. Hva kan tjene som en forklaring på at vi slik har «mistet» to observasjoner når vi skal estimere støyvariansen?

- a) Fordi vi har to variabler, Y og X som studeres.
b) Fordi støyvariansen beskriver variasjon rundt en rett linje, og det kreves minst to observasjoner for å plassere linja.
c) Fordi vi skal estimere både en forventning og et standardavvik.
d) Fordi SST kan splittes i to delkvadratsummer, SSR og SSE.
e) Fordi sammenhengen mellom X og Y enten er positiv eller negativ.

Emneansvarlig:

Solve Sæbø

Sensor:

Torfinn Torp



Kandidatnummer: _____

Svarskjema: Sett ett kryss i hver rad i den kolonnen som svarer til det alternativet du mener er riktig svar på spørsmålet. Det er kun tillatt å sette ett kryss i hver rad. (Dersom du vil endre svaret ditt, marker tydelig at du velger bort alternativet ved å skravere bort krysset.)

Oppgave	a	b	c	d	e
1					
2					
3					
4					
5					
6					
7					
8					
9					
10					
11					
12					
13					
14					
15					
16					
17					
18					
19					
20					
21					
22					
23					
24					
25					
26					
27					
28					
29					
30					
31					
32					
33					
34					
35					
36					
37					
38					
39					
40					

