

Obligatorisk øving for STK1110, Høsten 2020

Øving 2 av 2

Innleveringsfrist

Torsdag 5. november 2020, klokken 14:30 i Canvas (canvas.uio.no).

Instruksjoner

Du velger selv om du skriver besvarelsen for hånd og scanner besvarelsen eller om du skriver løsningen direkte inn på datamaskin (for eksempel ved bruk av LaTeX). Besvarelsen skal leveres som **én PDF-fil**. Scannede ark må være godt lesbare. Besvarelsen skal inneholde navn, emne og obliqnummer.

Det forventes at en har en klar og ryddig besvarelse med tydelige begrunnelser. Husk å inkludere alle relevante plott og figurer. Studenter som ikke får sin opprinnelige besvarelse godkjent, men som har gjort et reelt forsøk på å løse oppgavene, vil få én mulighet til å levere en revidert besvarelse. Samarbeid og alle slags hjelpemidler er tillatt, men den innleverte besvarelsen skal være skrevet av deg og reflektere din forståelse av stoffet. Er vi i tvil om du virkelig har forstått det du har levert inn, kan vi be deg om en muntlig redegjørelse. I oppgaver der du blir bedt om å programmere må du legge ved programkoden og levere den sammen med resten av besvarelsen.

Søknad om utsettelse av innleveringsfrist

Hvis du blir syk eller av andre grunner trenger å søke om utsettelse av innleveringsfristen, må du ta kontakt med studieadministrasjonen ved Matematisk institutt (e-post: studieinfo@math.uio.no) i god tid før innleveringsfristen. For å få adgang til avsluttende eksamen i dette emnet, må en bestå alle obligatoriske oppgaver i ett og samme semester.

Spesielt om dette oppgavesettet

Du skal bruke programpakken R til å gjøre beregninger i oppgavene, og du må angi hvilke kommandoer du har brukt for å komme fram til svarene dine. For å få godkjent besvarelsen, må du ha minst 65% riktig på hver av de 4 oppgavene.

For fullstendige retningslinjer for innlevering av obligatoriske oppgaver, se her:

www.uio.no/studier/admin/obligatoriske-aktiviteter/mn-math-oblig.html

LYKKE TIL!

Oppgave 1

Følgende tabell viser 10 målinger av kroppstemperatur for kvinner, x_1, \dots, x_{10} , og 10 målinger for menn, y_1, \dots, y_{10} . Hensikten med denne oppgaven er å undersøke om det er tilstrekkelig informasjon i tabellen til å kunne konkludere med at kroppstemperaturen er forskjellig for menn og kvinner. Dataene finnes som "temp.txt" i mappen <https://www.uio.no/studier/emner/matnat/math/STK1110/data/>.

| Kroppstemperatur | |
|------------------|---------|
| Menn | Kvinner |
| 36.1 | 36.6 |
| 36.3 | 36.7 |
| 36.4 | 36.8 |
| 36.6 | 36.8 |
| 36.6 | 36.7 |
| 36.7 | 37.0 |
| 36.7 | 37.1 |
| 37.0 | 37.3 |
| 36.5 | 36.9 |
| 37.1 | 37.4 |

- Lag et boksplokk som viser fordelingen av observasjonene. Kommentér hva du finner.
- Lag normalfordelingsplokk for de to observasjonssettene, altså ett for menn og ett for kvinner. Kommentér hva du ser.

I resten av oppgaven antar vi at observasjonene er realisasjoner av normalfordelte variabler. I c) og d) skal du forklare hvordan tester og konfidensintervaller konstrueres, og sette inn i formlene du utleder. Sjekk deretter svarene du får mot R-proseduren `t.test()`.

- Anta at variansen er den samme for de to utvalgene, og test med signifikansnivå 5% om det er noen forskjell i forventet kroppstemperatur. Beregn P-verdien, og lag et 95% konfidensintervall for denne forskjellen.
- Gjennomfør testen og beregn P-verdien også i det tilfellet der en ikke antar felles varians. Diskutér og forklar resultatene.
- Utleid og gjennomfør en F-test for å sjekke om det er noen grunn til å påstå at variansene er forskjellige. Sjekk mot `var.test()` i R.
- Se nå på situasjonen der en vurderer å innhente to nye målinger. La X_{11} være verdien for kvinnen og Y_{11} verdien for mannen, slik at forskjellen er $X_{11} - Y_{11}$. Vi antar nå at alle observasjonene er normalfordelte med samme varians. Begrunn at et rimelig anslag for $X_{11} - Y_{11}$ er differansen mellom gjennomsnittet av de 10 eksisterende målingene for kvinner og menn, altså $\bar{X} - \bar{Y}$. Hva er fordelingen til $X_{11} - Y_{11} - (\bar{X} - \bar{Y})$? Bruk dette til å lage et 95%

prediksjonsintervall for $X_{11} - Y_{11}$, altså et intervall som med sannsynlighet 0.95 inneholder $X_{11} - Y_{11}$. Dette er gjennomgått for ett-utvalgs-situasjonen på forelesning. Forklar hva som er forskjellen mellom et slikt intervall og et konfidensintervall for $\mu_1 - \mu_2$. Hvordan skal et prediksjonsintervall tolkes? [Hint: Siden alle variablene er normalfordelte, er $X_{11} - Y_{11} - (\bar{X} - \bar{Y})$ også det. Det er derfor nok å beregne forventning og varians for å finne fordelingen til denne størrelsen.]

Oppgave 2

Siden eneggede tvillinger har samme genetiske materiale, brukes såkalte tvillingstudier til å kartlegge hvordan miljøet virker inn på ulike egenskaper. I en bok av den amerikanske forskeren Susan Faber finner vi data for $n = 31$ tvillingpar, der den ene tvillingen vokste opp hos biologiske foreldre (Twin A) og den andre vokste opp hos andre familiemedlemmer, foster- eller adoptivforeldre (Twin B). Nedenfor finnes en oppsummering av målt IQ for disse personene. Spørsmålet vi ønsker å belyse er om det er forskjell i IQ hos eneggede tvillinger der den ene tvillingen har vokst opp hos biologiske foreldre, og den andre ikke.

| | N | Mean | StDev | SE Mean |
|------------|----|-------|-------|---------|
| Twin A | 31 | 93.32 | 15.41 | 2.77 |
| Twin B | 31 | 96.58 | 13.84 | 2.49 |
| Difference | 31 | -3.26 | 8.81 | 1.58 |

- Begrunn hvorfor en parett sammenligning er best egnet i denne situasjonen. Beskriv kort hvilke antakelser vi må legge til grunn for videre analyse.
- Kall forventet forskjell mellom Twin A og Twin B for μ_D . Sett opp nullhypotese og alternativ hypotese for å besvare spørsmålet om forskjell i IQ. Finn en egnet testobservator, og beregn dennes verdi. Beregn så tilhørende p-verdi. Spesifiser antall frihetsgrader i fordelingen du bruker. Formulér din konklusjon på testen.
- Finn et 95% konfidensintervall for μ_D . Hva betyr det at dette intervallet dekker kun negative verdier? Forklar kort om sammenhengen mellom tosidig testing og konfidensintervaller.

Oppgave 3

En undersøkelse presentert i Aftenposten slo opp på førstesiden at småbarnsfedrene nå opplever tidsklemma (mellom familie og arbeidsliv) sterkere enn småbarnsmødrene. Undersøkelsen bygde på intervjuer med 3000 kvinner og 3000 menn som har barn i rett alder. 16.2% av fedrene (dvs. 486 personer) opplevde ofte tidsklemmeproblemer, mens 14.7% (dvs. 441 personer) av mødrene opplevde det samme.

- Er forskjellen mellom mødre og fedre signifikant? Formulér hypoteser, beregn en p-verdi, og konkluder. Kommentér kort.

b) Kontrollér svaret ditt ved å bruke `prop.test()` i R.

Oppgave 4

Tabellen nedenfor angir 18 målinger av snømengde om vinteren i et fjellområde og vannstanden i en elv i samme område etter snøsmelting om våren. De 18 målingene representerer 18 sesonger spredd over en lengre tidsperiode. Her er vannstanden, som er angitt i tommer, respons- eller avhengig variabel, mens snømengden, målt ved noe som heter vannekvivalens, er forklaringsvariabel. Sammenhengen mellom snømengde og vannstand er viktig for bl.a. prediksjon av vannføring og flomfare. Dataene finnes i fila `snoe_vann.txt` på kursets hjemmeside.

| Snøinnhold | Vannstand |
|------------|-----------|
| 23.1 | 10.5 |
| 32.8 | 16.7 |
| 31.8 | 18.2 |
| 32.0 | 17.0 |
| 30.4 | 16.3 |
| 24.0 | 10.5 |
| 39.5 | 23.1 |
| 24.2 | 12.4 |
| 52.5 | 24.9 |
| 37.9 | 22.8 |
| 30.5 | 14.1 |
| 25.1 | 12.9 |
| 12.4 | 8.8 |
| 35.1 | 17.4 |
| 31.5 | 14.9 |
| 21.1 | 10.5 |
| 27.6 | 10.5 |
| 27.6 | 16.1 |

- Beskriv en enkel lineær regresjonsmodell for sammenhengen mellom snømengde og vannstand. Tilpass en regresjonslinje til dataene ovenfor ved hjelp av R-funksjonen `lm()`. Plott observasjonene og den tilpassede regresjonslinja. Virker estimatene for koeffisientene rimelige?
- Plott residualene mot forklaringsvariabelen. Lag også et normalfordelingsplott av residualene. Hvordan vurderer du modellens egnethet?
- Beregn et estimat for variansen til feilleddene. Konstruér et 95% konfidensintervall for stigningstallet β_1 .
- Utled en test for $H_0 : \beta_0 = 0$ mot $H_1 : \beta_0 \neq 0$ med signifikansnivå 5%. Gjennomfør testen. Hva er p-verdien? Sammenlign med resultatene fra `lm()`.