

Dette er oppgavesettet for andre obligatoriske innleveringsprosjekt for kurset STK 1110. Det legges ut på kursiden **torsdag 27/x/16**, og leveringsfristen er **torsdag 10/xi/16**. Besvarelsen skal leveres til instituttkontoret ved Matematisk institutt, senest kl. 14:30 den dagen. Sjekk de praktiske detaljene samlet i

[www.uio.no/studier/admin/obligatoriske-aktiviteter/mn-math-oblig.html](http://www.uio.no/studier/admin/obligatoriske-aktiviteter/mn-math-oblig.html),

der det også forklares at du skal bruke en bestemt standardisert «Forside for obligatoriske innleveringer»:

[www.uio.no/studier/admin/obligatoriske-aktiviteter/mn-math-obligforside.pdf](http://www.uio.no/studier/admin/obligatoriske-aktiviteter/mn-math-obligforside.pdf)

Hvis flere samarbeider om å løse oppgavene, må likevel hver student levere sin egen besvarelse. Det må gå frem av besvarelsen hvem du eventuelt har samarbeidet med.

Du kan levere håndskrevet eller maskinskrevet (tekstbehandlet) besvarelse, and you may write in bokmål, nynorsk, English or Latin. Der du bruker **R** må utskrifter og plott legges ved (eventuelt opereres inn i tekstbehandlingsdokumentet). Ulike introduksjonshefter for **R** er å finne på verdensveven; jfr. også det som gjennomgås i undervisningen.

**Nils Lid Hjort**

### Oppgave 1

FALLING IN LOVE IS LIKE ALGEBRA; you have to go step by step & follow the steps for it to turn out right\*. Resultater fra de noe algebraiske operasjonene vi skal gjennom her, skal om ikke annet hjelpe oss i oppgave 2. Jeg minner om at for en stokastisk variabel  $A$  med forventning  $\xi$ , så gjelder  $\text{Var } A = E(A - \xi)^2 = E A^2 - \xi^2$ . Spesielt er altså  $E A^2 = \xi^2 + \text{Var } A$ . Dessuten er det lov å vite, eller kjapt vise, uten at det er avkrevet her, at for vilkårlige tall  $x_1, \dots, x_n$  med gjennomsnitt  $\bar{x} = (1/n) \sum_{i=1}^n x_i$ , gjelder

$$\sum_{i=1}^n (x_i - \bar{x}) = 0 \quad \text{og} \quad \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - n\bar{x}^2.$$

- (a) La nå  $X_1, \dots, X_n$  være uavhengige stokastiske variable, med samme forventning  $\mu$  og samme varians  $\sigma^2$ . Vis at  $Z = \sum_{i=1}^n (X_i - \bar{X})^2$  har forventning  $(n-1)\sigma^2$ . Dette er grunnen at til man foretrekker  $n-1$  i nevneren for den vanlig definerte empiriske varians

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2,$$

for da blir altså  $E S^2 = \sigma^2$ .

---

\* or you just end up with a big headache, and wishing you hadn't done it in the start.

- (b) Anta nå at  $X_i$ -ene fortsatt har samme varians  $\sigma^2$ , men at de kan ha forskjellige forventninger, si  $\mu_1, \dots, \mu_n$ , med gjennomsnitt  $\bar{\mu} = (1/n) \sum_{i=1}^n \mu_i$ . Vis at da blir

$$E S^2 = \sigma^2 + \frac{1}{n-1} \sum_{i=1}^n (\mu_i - \bar{\mu})^2.$$

Den empiriske variansen griper altså fatt i ikke bare variablenes varians, per se, men i hvilken grad forventningene ikke er like.

- (c) Av og til må man arbeide med «nokså like» variable, men der disse ikke har helt like forventninger, og heller ikke helt like varianser. Anta derfor at  $X_i$ -ene er uavhengige, med forventninger  $E X_i = \mu_i$  og varianser  $\text{Var } X_i = \sigma_i^2$ . Vi skal arbeide med

$$Q = \sum_{i=1}^n \frac{(X_i - \hat{\mu})^2}{\sigma_i^2}, \quad \text{der} \quad \hat{\mu} = \frac{\sum_{i=1}^n X_i / \sigma_i^2}{\sum_{i=1}^n 1 / \sigma_i^2} = \frac{1}{A} \sum_{i=1}^n \frac{X_i}{\sigma_i^2},$$

der altså  $A = \sum_{i=1}^n 1 / \sigma_i^2$ . Vis at

$$E Q = n - 1 + \sum_{i=1}^n \frac{(\mu_i - \tilde{\mu})^2}{\sigma_i^2},$$

der

$$\tilde{\mu} = \frac{\sum_{i=1}^n \mu_i / \sigma_i^2}{\sum_{i=1}^n 1 / \sigma_i^2} = \frac{1}{A} \sum_{i=1}^n \frac{\mu_i}{\sigma_i^2}.$$

- (d) Skriv ut hva resultatet fra (c) gir, for de to tilfellene (i) alle  $\sigma_i$ -ene er like; (ii) alle  $\mu_i$ -ene er like.

## Oppgave 2

ET BILDE SIER MER ENN CA. FØRTI ORD, er noen villige til å mene, men hvor lange er disse ordene? Det vet studentene i dette kurs noe om, takket være oppgavesettet til Oblig I. Hver student skred i den forbindelse til sin bokhylle, fant 100 ord og deres lengder i en norsk bok, og 100 ord og deres lengder i en engelsk bok. Vi har ikke tatt vare på absolutt all informasjonen i disse  $2 \times 64 \times (100 + 100)$  ord, men på essensen, nemlig gjennomsnitt  $\bar{x}$  og empirisk standardavvik  $\hat{\sigma}$ , for hver av  $n = 64$  studenter, for de to språk. Datafilen `wordlengths-data` på kurssiden inneholder

$$(\bar{x}_{i,N}, \hat{\sigma}_{i,N}, \bar{x}_{i,E}, \hat{\sigma}_{i,E}) \quad \text{for } i = 1, \dots, n.$$

Om man lagrer denne datafilen på eget område, skal man kunne aksessere disse dataene ved å skrive

```
wordlengths <- matrix(scan("wordlengths-data", skip=10), byrow=T, ncol=5)
```

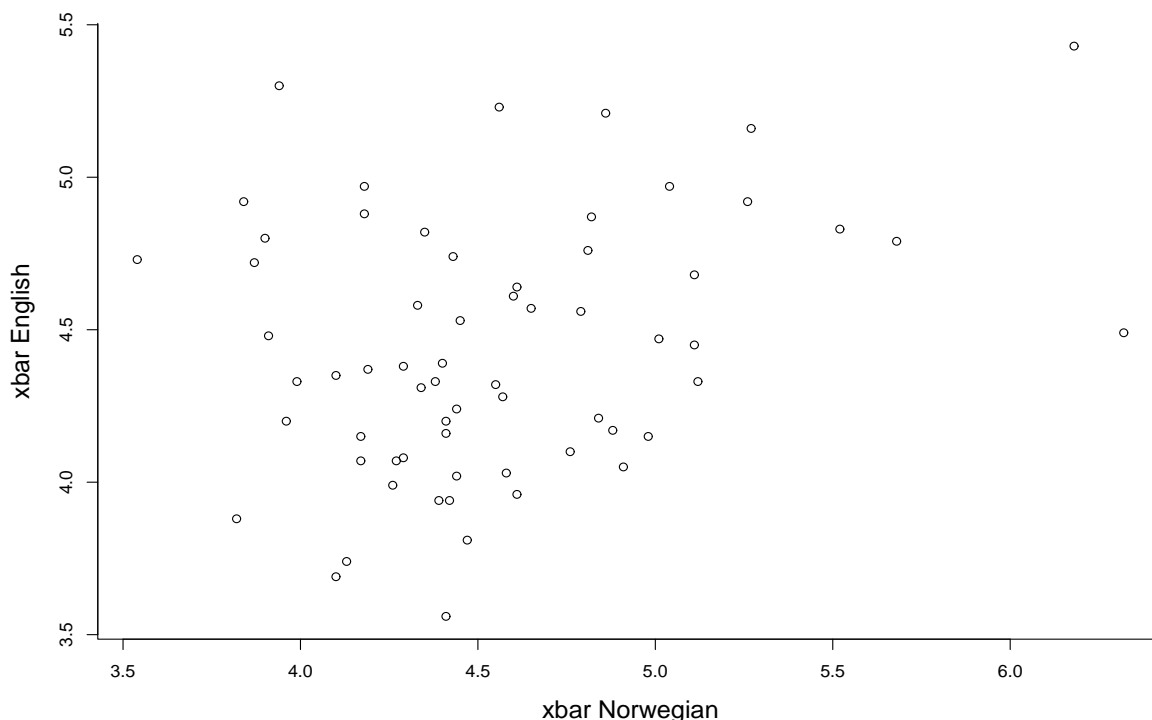
Etter dette kan man f.eks. skrive

```

xbarN <- wordlengths[ ,2]
sigmaN <- wordlengths[ ,3]
xbarE <- wordlengths[ ,4]
sigmaE <- wordlengths[ ,5]

```

fulgt av `plot(xbarN,xbarE)` og diverse statistiske operasjoner; se figur 1, der du til og med skal kunne se ditt eget bidrag til statistisklingvistisk erkjennelse.



Figur 1:  $(\bar{x}_{i,N}, \bar{x}_{i,E})$ , gjennomsnittssordlengder for norske og engelske bøker, for  $n = 64$  lesende studenter.

- Stoltenberg og Hjort har altså tillatt seg å ikke ta vare på mer enn nettopp  $(\bar{x}, \hat{\sigma})$ , for hver av studentene, for hvert av de to språk. I hvilken forstand kan de unne seg å ha god statistisk samvittighet, i denne anledning?
- Lag et slikt plott, over  $(\bar{x}_{i,N}, \bar{x}_{i,E})$ . Gi argumenter som gjør det rimelig å anta, ihvertfall som en approksimasjon, at disse parene kommer fra en binormal fordeling. Denne har en viss korrelasjonskoeffisient,  $\rho$ . Estimer denne, med den vanlige empiriske korrelasjonskoeffisienten,  $\hat{\rho}$ . Du skal også lage et tilnærmet 95% konfidensintervall for  $\rho$ , ved å anvende «transformasjonsknepet», at

$$A(\hat{\rho}) \approx_d N(A(\rho), 1/(n-3)), \quad \text{der} \quad A(\rho) = \frac{1}{2} \log \frac{1+\rho}{1-\rho}.$$

Gi en tolkning av det du kommer frem til.

- Er gjennomsnittlig ordlengde omtrent lik, for norsk og engelsk? Undersøk dette, og formuler en konklusjon.

- (d) Lag også et plott over standardavvikene, altså  $(\hat{\sigma}_{i,N}, \hat{\sigma}_{i,E})$ . For å belyse om gjennomsnittlig variansparameter for norske og engelske ordlengder er like, kan du her anvende at hvert estimert standardavvik,  $\hat{\sigma}_{i,N}$  og  $\hat{\sigma}_{i,E}$ , har en tilnærmet normalfordeling. (Du kan altså anvende metoder for tester og konfidensintervaller som tar utgangspunkt i dette, uten at du trenger å komme inn på de eksakte fordelinger for de  $2 \times n$  empiriske standardavvikene.) Som for (c), formuler en konklusjon basert på det du finner ut her.
- (e) Men er studentene, det vil si deres bokhyller, omtrent like, eller er det systematiske forskjeller mellom dem? Foretrekker noen studenter knappordet poesi, der andre best liker langordsfrekvente middelalderromaner? – En måte å nærme seg dette spørsmålet på er som følger. Vi starter med

$$\bar{x}_i \sim N(\mu_i, \kappa_i^2) \quad \text{for } i = 1, \dots, n$$

(for eksempel for det norske utvalget, med helt tilsvarende for det engelske), der vi også vet at  $\kappa_i^2 = \sigma_i^2/100$ , med gode estimater for  $\sigma_i$ -ene. Hvilken tolkning har  $\mu_i$  her? Så tenker vi oss at  $\mu_i$ -ene stammer fra sin egen normalfordeling,  $\mu_i \sim N(\mu_0, \tau^2)$ . Spørsmålet er om  $\tau$  er null eller om ikke annet temmelig liten, eller om den er merkbar. Hvilken tolkning vil du gi  $\tau$ ?

- (i) Anta i første omgang at  $\sigma_i$ -ene er tilnærmet like, og beregn

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (\bar{x}_i - \bar{x})^2,$$

der  $\bar{x}$  er det store gjennomsnittet av gjennomsnittene. Bruk så resultater fra oppgave 1 til å estimere  $\tau$ , for det norske og det engelske ordutvalget. Gi også en test for hypotesen om full likhet mellom  $\mu_i$ -ene (som svarer til at  $\tau = 0$ ).

- (ii) Prøv så å gjennomføre en analog, men litt mer kompleks analyse, igjen ved å bruke resultater fra oppgave 1, eller noe analogt, men der du ikke antar at  $\sigma_i$ -ene er like.

### Oppgave 3

«BRA PÅVIST! PÅVIS MER!», skrev altså *Natt og Dag* begeistret, for noen uker siden, om de temaene min PhD-student Céline Cunen har arbeidet med, der hun sammenligner aspekter ved rosekrigene (1455–1487) og *Game of Thrones* (det litterære univers skapt av George R.R. Martin). Dette har vi arbeidet med i Oblig I, og i denne oppgaven skal skal vi fortsette, i to nye retninger. Den første er å anvende gamma-fordelingen, med momentbaserte estimatorer (der Oblig I kjørte Weibull, med kvantilbasert estimering); den andre er bootstrapping. Se også Cunens bloggpost i FocuStat-prosjektet, og kom på hennes faglig-pedagogiske foredrag 3/11.

Gamma-fordelingen har tetthet

$$g(x, a, b) = \frac{b^a}{\Gamma(a)} x^{a-1} \exp(-bx) \quad \text{for } x > 0. \quad (\spadesuit)$$

Her er  $(a, b)$  positive parametre, og  $\Gamma(a)$  er den såkalte gamma-funksjonen. Læreboken anvender en noe annerledes og mindre vanlig parametrisering, men jeg velger altså den i (♠), og som også er den som er benyttet i **R** via `dgamma`, `pgamma`, `qgamma`, `rgamma`.

(a) Vis at forventning og varians for denne fordelingen blir

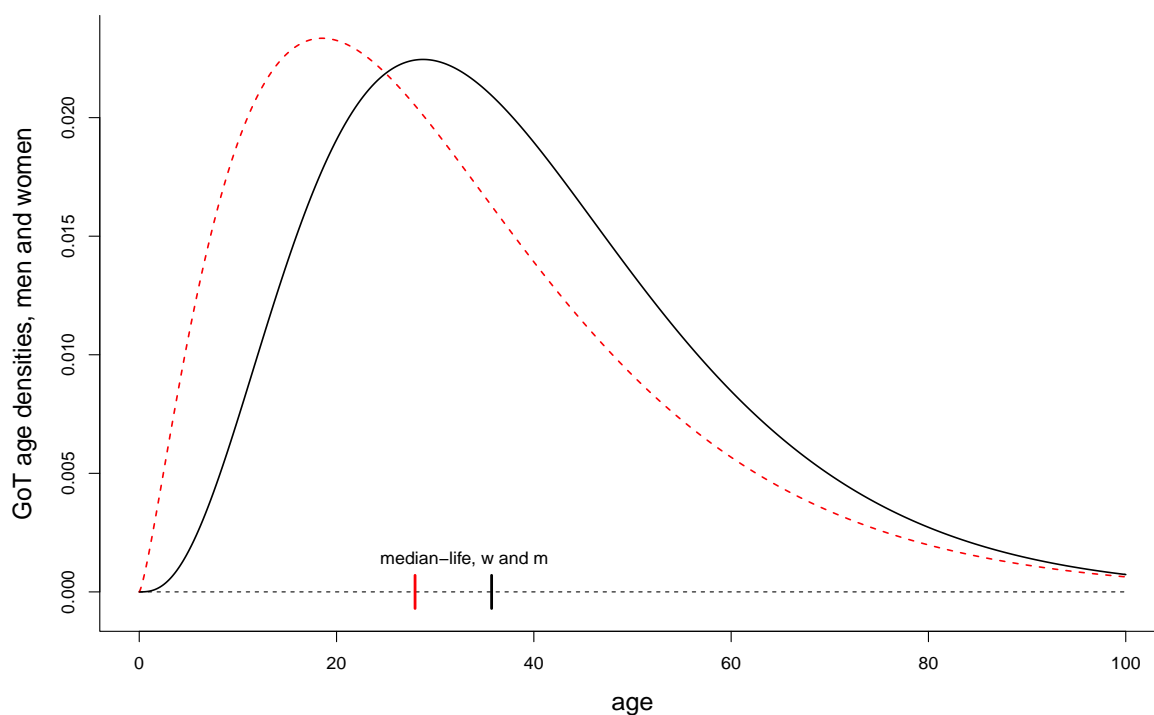
$$E X = \frac{a}{b} \quad \text{og} \quad \text{Var } X = \frac{a}{b^2}.$$

Her kan du få bruk at  $\Gamma(u+1) = u\Gamma(u)$ , for hver  $u$ , slik at også  $\Gamma(u+2) = (u+1)u\Gamma(u)$ , osv.

(b) Anta  $X_1, \dots, X_n$  er uavhengige observasjoner fra en slik gamma-fordeling. Forklar at momentmetoden for å estimere  $(a, b)$  består i å løse de to ligningene

$$\bar{X} = \frac{\hat{a}}{\hat{b}} \quad \text{og} \quad S^2 = \frac{\hat{a}}{\hat{b}^2},$$

der som vanlig  $\bar{X}$  er gjennomsnittet og  $S^2$  den empiriske varians. Løs disse, og sett opp eksplisitte formler for estimatorene  $(\hat{a}, \hat{b})$ .



Figur 2: Estimerte tettheter  $\hat{f}_m(x) = g(x, \hat{a}_m, \hat{b}_m)$  for menn og  $\hat{f}_w(x) = g(x, \hat{a}_w, \hat{b}_w)$  for kvinner, i GoT-universet, via gamma-fordelinger. De assosierte medianlevetidene er også tegnet inn.

(c) Få igjen tak i levetidsdataene for de 26 kvinner og de 127 menn, fra GoT-universet, via oppskriften fra Oblig I. Tilpass gamma-fordelingen til hver av disse datasettene, ved å finne momentestimatene  $(\hat{a}_m, \hat{b}_m)$  for menn og  $(\hat{a}_w, \hat{b}_w)$  for kvinner. Lag en versjon av figur 2, og kommenter kort hva du finner ut.

- (d) Beregn de estimater for medianlevetiden,  $q_w$  for kvinner og  $q_m$  for menn, som følger av de estimerte gamma-fordelinger. Det finnes ingen eksplisitt formel for

$$\hat{q}_m = G^{-1}\left(\frac{1}{2}, \hat{a}_m, \hat{b}_m\right),$$

der  $G(x, a, b)$  er den kumulative gamma-fordelingen, men det er ingen sak for **R**, der vi kan bruke

```
qhatw = qgamma(0.50,ahatw,bhatw)
```

```
qhatm = qgamma(0.50,ahatm,bhatm)
```

Disse er tegnet inn i figur 2.

- (e) Så skal vi bruke bootstrapping for å lage konfidensintervaller for  $q_w$  og  $q_m$ . Dette kan foregå omtrent slik, i **R**, der maskinen først har fått vite om datasettet `xmen` for mennene, med sample size `nmen` (lik 127):

```
boot <- 1000
```

```
qhatmstar <- 0*(1:boot)
```

```
for (b in 1:boot)
```

```
{
```

```
xmenstar <- sample(xmen,nmen,replace=T)
```

```
# this is a bootstrap sample of the xxm dataset
```

```
[then a few lines and formulas here
```

```
to compute ahatmstar and bhatmstar based on dataset xmenstar
```

```
using the same formulas for parameter estimates as in (b)]
```

```
#
```

```
qhatmstar[b] <- qgamma(0.50,ahatmstar,bhatmstar)
```

```
}
```

Da har du fått laget deg `boot` bootstrapverdier  $\hat{q}_m^*$ , og da gjenstår det bare å be om `quantile(qhatmstar,c(0.05,0.95))`

Bruk dette oppsettet til å finne tilnærmede 90% konfidensintervaller for medianlevetidene  $q_w$  og  $q_m$ . Kommenter det du finner ut.

- (f) Det er verd å registrere at bootstrapmetoden her er temmelig generell, og temmelig uavhengig av matematiske formler. Altså kan man anvende metoden også for mer komplekse problemstillinger enn dette. Velg *en statistisk parameter til*, etter eget ønske, og bruk oppsettet over til å finne estimater og konfidensintervaller, for din parameter, for GoT-populasjonene kvinner og menn. Kommenter igjen det du finner frem til.