

## Øving 10, ST1301

**Oppgave 1** Øyvind går to skiturer i påsken. På den ene skituren viste klokken 12:32 når han startet og 12:59 når han kom tilbake på hytta. Den andre dagen viste klokken 12:32 når han startet og 12:58 når han vendte tilbake. Klokken hennes viser bare timer og minutter og ikke sekunder. Hva er sannsynligheten for at skituren den andre dagen var lenger (målt i tid) enn skituren den første dagen?

Hint: Når klokken viser f.eks. 12:32 betyr det at tiden  $T_1$  var et sted mellom 12:32:00 og 12:33:00 ved avlesning av starttidspunkt første dag. Bruk uniforme sannsynlighetsfordelinger for å modellere hva vi vet om  $T_1$  og de tre andre tidspunktene som inngår i problemet. La f.eks.  $T_1$  være uniform på intervallet fra 32 til 33.

Finn sannsynligheten for hendelsen over (lengst skitur andre dag) ved å simulere 10000 realisasjoner av  $T_1, T_2, \dots, T_4$  og tell opp antall ganger hendelsen inntreffer f.eks. ved bruk av tellevariabel og if-setning.

**Oppgave 2** *Programmering av likelihood funksjon. Beregning av SME.*

I denne oppgaven skal vi estimere når og hvor lange brunstperioden er i en elgpopulasjon. Dette vil naturlig nok ikke være direkte observerbart men kan estimeres fra relevante data. Her skal vi bruke data gjort tilgjengelige av Erling Solberg ved NINA i Trondheim. Last ned dataene `ovul2.dat` fra hjemmesiden, og last disse inn i R ved å bruke `read.table` og `attach`. Disse dataene består av tre variable.  $t_1, t_2, \dots, t_n$  (har navnet `tid` i data.framen) er antall dager siden nyttår jaktdag  $i$  fant sted. Variablen  $n_i$  er antall elgkyr feltt på dag  $i$ , og variabelen  $x_i$  er antall ut av disse som hadde ovulert (hatt egglosning) på fellingstidspunktet. (Dette ble bestemt i laboratoriet ved å undersøke om en spesiell struktur (en eller flere corpus luteum.<sup>1</sup> var tilstede i ovariene.). Hos elg inntreffer ovulasjon en (eller flere ganger hvis først ovulasjon ikke fører til befruktning) i løpet av høsten for kjønnsmodne individer. Dermed vil andelen individer som har ovulert øke utover høsten og forløpet på denne sammenhengen vil kunne fortelle oss om gjennomsnittet og standardavviket til fordelingen av ovulasjonstidspunktene i populasjonen som vi her ønsker å estimere.

Lag først et spredningsplot av andelen ( $x_i/n_i$ ) som har ovulert på ulike tidspunkt  $t_i$ . La størrelsen på hvert punkt i plottet avhenge av  $n_i$  slik at det fremkommer av plottet hvor mange observasjoner hvert punkt representer. Dette kan gjøres via argumentet `cex`. (Se `?plot.default`.)

La oss nå anta at ovulasjonstidspunktet  $T$  til et gitt individ i populasjonen er normalfordelt med forventning  $\mu$  og standardavvik  $\sigma$ . Ikke alle kyrne

---

<sup>1</sup>Se [http://en.wikipedia.org/wiki/Corpus\\_luteum](http://en.wikipedia.org/wiki/Corpus_luteum).

ovulerer fordi ikke alle har nådd kjønnsmodning som 3-åringer. La oss derfor anta at hvert individ ovulerer en gang i løpet av høsten med sannsynlighet  $q$ . Disse antakelsene utgjør tilsammen vår modell. Vi ønsker nå å sette opp likelihoodfunksjonen for de observerte dataene for denne modellen. Vi trenger altså å finne et uttrykk for sannsynligheten for dataene gitt parameterne  $\mu$ ,  $\sigma$ , og  $q$ .

La  $O$  være hendelsen at et gitt individ ovulerer en eller annen gang i løpet av høsten. Sannsynligheten for at ovulasjon har funnet sted ved tidpunkt  $t_i$ , altså at  $T < t_i$  blir da

$$\begin{aligned} P(T \leq t_i | O_i) &= P\left(\frac{T - \mu}{\sigma} \leq \frac{t_i - \mu}{\sigma}\right) \\ &= P\left(Z \leq \frac{t_i - \mu}{\sigma}\right) \\ &= G\left(\frac{t_i - \mu}{\sigma}\right) \end{aligned} \tag{1}$$

hvor  $G$  er kumulativ tetthet til en standard normalfordelt variabel.

I likelihoodfunksjonen vår trenger vi de ubetingede sannsynligheten for hendelsen  $T < t_i$ . La oss kalle denne sannsynligheten for  $p_i$ . Bruker vi lov om totalsannsynlighet får vi at

$$\begin{aligned} p_i &= P(T < t_i) \\ &= P(T < t_i | O_i)P(O_i) + P(T < t_i | \bar{O}_i)P(\bar{O}_i) \\ &= G\left(\frac{t_i - \mu}{\sigma}\right)q + 0 \cdot (1 - q) \end{aligned} \tag{2}$$

i og med at hendelsen  $T < t_i$  har sannsynlighet 0 hvis individ  $i$  ikke ovulerer på noe tidspunkt i løpet av høsten (hendelsen  $\bar{O}_i$ ).

Ser vi på antall individer som har ovulert  $X_i$  på et bestemt dag  $i$ , hver med sannsynlighet  $p_i$ , for vi at  $X_i$  må være binomisk fordelt med parameter  $p_i$  og  $n_i$ . Dermed blir likelihoodfunksjonen (sannsynligheten for alle de observerte dataene)

$$L(\mu, \sigma, q) = \prod_{i=1}^n \binom{n_i}{x_i} p_i^{x_i} (1 - p_i)^{n_i - x_i}, \tag{3}$$

hvor hver  $p_i$  er en funksjon av  $\mu$ ,  $\sigma$  og  $q$  gitt ved ved ligning (2).

Lag et plot av  $p_i$  som funksjon av  $t_i$  for  $\mu = 260$ ,  $\sigma = 10$  og  $q = 0.8$  (altså funksjonen gitt ved (2)). Hvordan passer funksjonen med spredningsplottet av  $x_i/n_i$  og  $t_i$ ?

Programmer så en funksjon som beregner verdien av minus log likelihoodfunksjonen for gitte parameterverdier. For å fungere sammen med `optim` må

likelihoodfunksjonen må ta i mot parameterverdiene som første argument i form av en vektor, samt vektorer  $\mathbf{t}$  og  $\mathbf{x}$  som inneholder de observerte dataene. Tips: Beregn  $p_i$ 'ene og ta vare på dette som et mellomregningsresultatet (bruk `pnorm`). Beregn så log-likelihoodverdien ved bruk av elementvise operasjoner på vektorer, funksjonen `sum`, og eventuelt `dbinom` med tilleggsargumentet `log=T`.

Beregn så sannsynlighetsmaksimeringsestimaterne av  $\mu$ ,  $\sigma$ , og  $q$  ved bruk av `optim`. Det kan være viktig å velge gode startverdier for parameterne. Kontroller at `optim` finner samme maksimum om du endrer på startverdiene.

Hva blir SMEene av parameterne i modellen? Legg inn et plot av funksjonen gitt ved (??) i spredningsplottet for  $t_i$  og  $x_i$  ved å bruke `lines`. Hvordan passer den tilpassede modellen med dataene?

Dersom vi betrakter et tidintervall slik at 95% av individene i populasjon ovulerer i dette intervallet, hvor langt må da intervallet være?

**Oppgave 3** Oppgave 5.4 i ISwR. Se eventuelt også notatet om binormalfordelingen fra i høst.