

Obligatorisk øvelse 3 i STAT110/622

Basert på kapittel 1-8 i boken. Les informasjon “Obligatoriske innleveringer” på MittUiB->Hjem.

1. Kollevåundersøkelsen

Start med å se videoen som det er linket til på MittUiB. Det ble tatt prøver på fire ulike steder (stasjoner) som vi kaller 1, 2, 3 og Referanse (Ref). Referanseprøvene ble tatt langt borte fra Kollevåg på et sted med minimal forurensing. Figur 1 viser empirisk gjennomsnitt og standardavvik for ensymnivået for et enzym som forkortes Cat. (Merk at figuren hverken er et histogram eller et boxplot.) Spørsmålet vi er interessert i er om ensymnivået inne i Kollevåg (stasjon 1–3) er ulikt det referansestasjon (Ref).

(a) Les av \bar{x} og s fra Figur 1 for hver av stasjonene 1, 2 og 3.

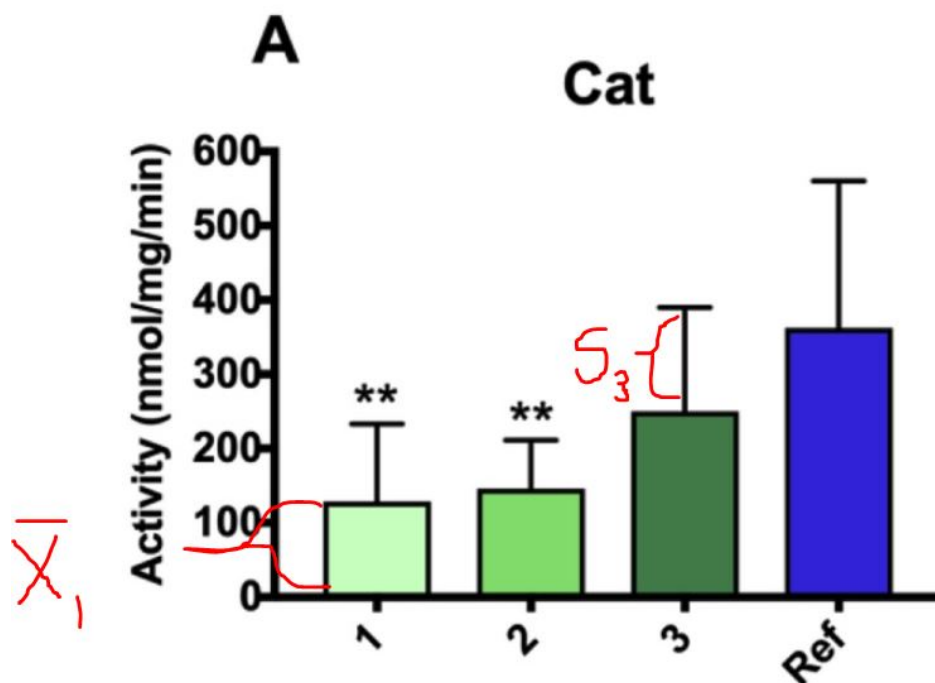


Figure 1: Gjennomsnitt (\bar{x}) og standardavvik (s) for prøver tatt fra fire stasjoner i Kollevåg. (Det er tegnet inn for hånd hvordan du leser av \bar{x} og s .) Antall fisk som det er tatt prøver fra er $n_1 = 18$, $n_2 = 20$, $n_3 = 13$ og $n_{\text{Ref}} = 17$. Data er hentet fra den vitenskapelige publikasjonen Dale et al. “Contaminant accumulation and biological responses in Atlantic cod (*Gadus morhua*) caged at a capped waste disposal site in Kollevåg, Western Norway.” Marine environmental research 145 (2019): 39-51.

Vi antar at variasjonen mellom individer innen en stasjon følger en normalfordeling, og vi betegner forventingsverdien for hver av de fire stasjonene med μ_1 , μ_2 , μ_3 og μ_{Ref} . Den praktiske fortolkningen av en μ er gjennomsnittlig ensymnivå for alle individer som “lever” på den stasjonen.

- (b) Lag et 95% konfidensintervall for hver av μ_1 , μ_2 og μ_3 (totalt tre intervaller).

Vi leser av $\bar{x}_{\text{Ref}} = 354$ fra Figur 1, og vi skal nå ignorere usikkerheten i denne målingen ved å sette $\mu_{\text{Ref}} = 354$, dvs. anta at dette er gjennomsnittet blant alle individer som lever på referansestasjonen.

- (c) Hvilke av de tre konfidensintervallene som du konstruerte under punkt b) inneholder $\mu_{\text{Ref}} = 354$?
Hvis du har lest kapittel 9: gjenkjenner du dette spørsmålet som en hypotesetest (hvilken hypotese)?

2. Sannsynlighetsmaksimering

En postfunksjonær har en betjeningstid T som er eksponensialfordelt med tetthet

$$f(t) = \theta e^{-\theta t}, \quad t \geq 0,$$

der θ er en ukjent parameter.

- (a) Gitt n observasjoner t_1, \dots, t_n , finn sannsynlighetsmaksimeringsestimatet $\hat{\theta}$ for θ . Finn en numerisk verdi for $\hat{\theta}$ når en har følgende 10 observerte betjeningstider:

t_i : 1, 1.4, 2.0, 0.5, 0.7, 2.0, 1.3, 1.1, 1.8, 0.2,

der tidsenhet er ett minutt.

- (b) Er $\hat{\theta}$ en forventningsrett estimator for θ ? Begrunn svaret.

3. Sannsynlighetsmaksimering

Vi har uavhengige tilfeldige variable X og Y med fordeling $X \sim \text{Poisson}(\theta)$ og $Y \sim \text{Poisson}(2\theta)$, og observasjoner $x = 3$ og $y = 5$ av disse.

- (a) Vis at uttrykket for log-likelihood funksjonen er gitt ved

$$l(\theta) = [5 \ln(2) - \ln(3!) - \ln(5!)] + 8 \ln \theta - 3\theta.$$

Plott $l(\theta)$ for verdier av $\theta \in [0, 10]$ (en grov skisse er OK). For ca hvilken verdi av θ har $l(\theta)$ sitt maksimum?

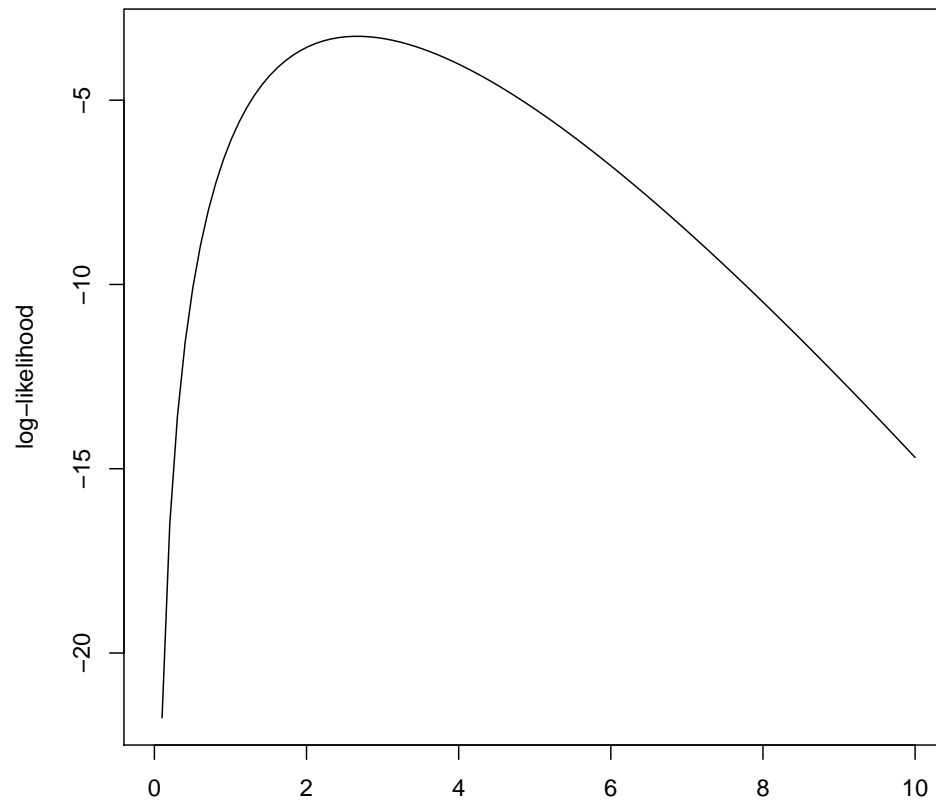


Figure 2: Log-likelihood funksjonen plottet som funksjon av θ

- (b) Beregn sannsynlighetsmaksimeringsestimatoren $\hat{\theta}$ for de observerte verdiene $x = 3$ og $y = 5$.

4. **Hvilken estimator er best?** Et politisk parti finansierer sin virksomhet ved hjelp av et lotteri. Lotteriet er organisert slik at hver deltager har mulighet til å få gevinst i tre uavhengige spilleomganger. I følge reklamemateriellet fra partiet er det ulik sjanse for gevinst i de tre omgangene. Hvis vi betegner sannsynligheten for gevinst i første omgang med θ , er sannsynligheten for gevinst i andre omgang lik 4θ , og i tredje omgang 5θ . Her er θ en parameter som tilfredsstiller $\theta < \max(1/4, 1/5)$. Både i andre og tredje omgang har alle deltagerne de samme gevinstmulighetene, uansett om vedkommende allerede har vunnet gevinst i en tidligere omgang eller ikke.

Vi innfører de tre tilfeldige variablene X_1, X_2, X_3 som angir om du får gevinst i hver av de tre omgangene. Vi lar X_j være lik 1 dersom j 'te omgang gir gevinst, og 0 ellers, for $j = 1, 2, 3$.

Hint: Oppgaven likner på oppgave 3 vår eksamen 2017, og du kan bruke denne som utgangspunkt for løsning om du vil.

- (a) Vis at disse har forventning og varians:

$$\begin{aligned} E(X_1) &= \theta, & V(X_1) &= \theta(1 - \theta) \\ E(X_2) &= 4\theta, & V(X_2) &= 4\theta(1 - 4\theta) \\ E(X_3) &= 5\theta, & V(X_3) &= 5\theta(1 - 5\theta) \end{aligned}$$

Du vil bestemme et estimat for sannsynligheten θ ut fra de tre observasjonene X_1, X_2, X_3 . Du foreslår at man skal bruke en av disse estimatorene:

$$\hat{\theta}_1 = \frac{1}{10} (X_1 + X_2 + X_3), \quad \hat{\theta}_2 = \frac{1}{3} \left(X_1 + \frac{X_2}{4} + \frac{X_3}{5} \right).$$

- (b) Vis at begge estimatorene er forventningsrette. Finn uttrykk for variansene til $\hat{\theta}_1$ og $\hat{\theta}_2$ uttrykt ved θ . Avgjør hvilken estimator som er best når $\theta = 0.2$.

5. Gartneri

Ved et gartneri dyrkes hodekål. Vekten av den modne kålen, X , antas å være normalfordelt med forventningsverdi $\mu = 2.2$ kg og standardavvik $\sigma = 0.8$.

- (a) Hvis en velger et kålhode tilfeldig, hva er sannsynligheten for at dette skal: i) veie mindre enn 1.5 kg? ii) veie mellom 2 og 2.5 kg? Hva er sannsynligheten for at vektforskjellen mellom to tilfeldig valgte kålhoder skal være mer enn 1 kg?

Kålhoder som veier mindre enn 1.5 kg oppfyller ikke kravet til klasse l-kål.

- (b) Gitt at et kålhode oppfyller kravet (veier minst 1.5 kg), hva er sannsynligheten for at det veier mellom 2 og 2.5 kg?

Gartneriet bestemmer seg for å prøve ut en ny kåltype. Det påstas at hodene av denne nye typen gjennomgående veier mer enn hodene av den gamle typen. En liten prøveavling med 10 planter dyrkes. Vi antar først at de 10 vektmålingene Y_1, Y_2, \dots, Y_{10} er uavhengige og normalfordelte med ukjent forventning μ_Y og kjent standardavvik $\sigma_Y = 0.8$.

Resultatet av målingene er gitt i tabellen under:

kål nr. (i)	1	2	3	4	5	6	7	8	9	10
y_i	3.7	2.2	1.8	4.5	2.4	2.5	2.6	2.1	2.5	2.2

- (c) Hvilken estimator ville du bruke for å estimere μ_Y ? Skriv opp estimatoren og regn ut estimatet fra dataene i tabellen. Finn et 90%-konfidensintervall for μ_Y . Hva blir lengden av konfidensintervallet? Hvor mange planter måtte vi minst hatt for at lengden av konfidensintervallet skulle blitt mindre enn 0.2 kg?

6. Målinger

I en kommune finnes det en målestasjon for luftforurensing. Fra denne får vi i løpet av en dag 5 målinger: Y_1, Y_2, \dots, Y_n som antas uavhengige og identisk normalfordelte med forventning μ og varians σ^2 , der både μ og σ er ukjente parametre. Her vil μ angi graden av forurensning. Anta at $n = 5$ og at det en dag gjøres følgende målinger: 252, 311, 268, 287, 302. For senere bruk oppgis at $\sum_{i=1}^5 (y_i - \bar{y})^2 = 2342$.

- (a) Finn et 95% konfidensintervall for μ basert på dataene er som gitt ovenfor. Hvordan fortolker du en dekningsgrad på 95%?