

# STK1000: Løsningsforslag Uke 46

2022

## Check In - Oppgave 11.1

- Responsvariabelen er endelig eksamensresultat (final exam score).
- Case*'ene er observasjonene  $(y_i, x_{i1}, x_{i2}, \dots)$  altså kombinasjonen av respons og forklaringsvariabler for hver enkelt observasjon. Vi har totalt  $n = 166$  observasjoner.
- Vi har  $p = 7$  forklaringsvariabler.
- Vi har forklaringsvariablene "math course anxiety", "math test anxiety", "numerical task anxiety", "enjoyment", "self-confidence", "motivation" og "perceived usefulness of the feedback sessions".

## Check In - Oppgave 11.2

- Vi har  $\hat{y} = -10.8 + 3.2 \cdot 4 + 2.8 \cdot 2 = 7.6$ .
- Nei, det holder at forklaringsvariablene ligger i nærheten av observasjonene i data settet får å gi en fornuftig prediksjon  $\hat{y}$ .
- For en fiksert  $x_1$  vil  $\hat{y}$  forandre seg med 2.8 enheter for hver enhet  $x_2$  forandrer seg. Så for en forandring i  $x_2$  på 3 enheter, vil  $\hat{y}$  forandre seg med  $2.8 \cdot 3 = 8.4$  enheter.

## Oppgave 14.11

Kapittel 14 er ikke i den trykte utgaven, men en online versjon er tilgjengelig her <http://bcs.whfreeman.com/webpub/statistics/ips9e/9781319013387/companionchapters/companionchapter14.pdf>

- Proporsjonen er  $\hat{p} = 462/1003 = 0.46$ .
- Vi har  $\text{odds} = \frac{\hat{p}}{1-\hat{p}} = \frac{0.4606}{1-0.4606} = 0.8539$ .

## Oppgave 14.13

- Modellen er  $\log(\text{odds}_i) = \log \frac{p_i}{1-p_i} = \beta_0 + \beta_1 x_i$ . Her er  $x_i = 1$  hvis alder er mindre eller lik 25,  $x_i = 0$  hvis alder er større enn 25, og  $p_i$  er sannsynligheten for at individet har brukt mobiltelefonen i en butikk ['within the last 30 days to call a friend or family member for advice about a purchase they were considering'], og  $\text{odds}_i = \frac{p_i}{1-p_i}$ .
- $\beta_0$  er log-oddsen for at personer over 25 år har brukt mobilen i en butikk ['within the last 30 days to call a friend or family member for advice about a purchase they were considering'], mens  $\beta_1$  er differansen i log-odds for personer under 25 år sammenlignet med personer over 25 år. En alternativ måte å forklare rollen til  $\beta_1$  er at  $\frac{\text{odds}_{x=1}}{\text{odds}_{x=0}} = \exp(\beta_1)$ , det vil si at personer under 25 år har odds som er  $\exp(\beta_1)$  ganger oddsene til personer over 25.

## Oppgave 14.14

- Vi har nå modellen  $\log(\text{odds}_i) = \log \frac{p_i}{1-p_i} = \beta_0 + \beta_1 x_i$ , hvor  $x_i$  alder og  $p_i$  er sannsynligheten for at individet har brukt mobiltelefonen i en butikk ['within the last 30 days to call a friend or family member for advice about a purchase they were considering'].

- b)  $\beta_1$  er nå gjennomsnittlig forandring i log-odds for hvert år. Det vil si at oddsen for å bruke mobilen i en butikk [‘within the last 30 days to call a friend or family member for advice about a purchase they were considering’] multipliseres med en faktor  $\exp(\beta_1)$  per økning i alder på et år.
- c) Denne modellen antar at det er et lineært forhold mellom alder og log-odds. Det var ikke tilfellet i Oppgave 14.13 etter som vi bare hadde en enkelt indikator variabel for å skille mellom to aldersgrupper. For å undersøke denne antagelsen kan vi lage et plott tilsvarende det i Example 14.8, med alder langs x-aksen, og se om den estimerte linjen følger punktene. Merk at et slikt plott krever at vi har flere observasjoner for hver alder og at ikke alle observasjonene for en alder har samme respons, ettersom det vil gi en log-odds på  $\pm\infty$ .

## Eksamenoppgaver: 2009 oppg. 2, 2015 oppg. 3, 2016 oppg. 3 (alle disse er fra kap 11)

Løsningsforslag til eksamensoppgavene er på emnets semesterside <https://www.uio.no/studier/emner/matnat/math/STK1000/oppgaver/losningsforslag/>

### I tillegg:

**Oppgave:** Vis at utvalgs-gjennomsnittet  $\bar{x}$  faktisk er en maksimum likelihood-estimator for den sanne men ukjente parameteren forventningsverdien  $\mu$  i en  $N(\mu, \sigma)$  fordeling, basert på  $n$  uavhengige, identisk fordelte observasjoner  $x_1, x_2, \dots, x_n$ . NB! Da trenger du formelen for tetthetsfunksjonen, som du blant annet finner på side 50 i læreboka. Anta for enkelthets skyld at verdien til  $\sigma$  er kjent.

### Løsning:

La  $\mathbf{x}$  representere  $x_1, x_1, \dots, x_n$  (for litt mer kompakt notasjon). Vi har at tetthetsfunksjonen er identisk for alle observasjonene siden de har samme (identisk) fordeling, og at likelihooden er produktet av tetthetsfunksjonene siden observasjonene er uavhengige. Det er ofte lettere å maksimere log-likelihooden enn å maksimere likelihooden, og siden disse to er maksimale i samme verdi for  $\mu$ , så setter vi opp uttrykket for log-likelihooden før vi tenker på å optimere (finne maksimum):

$$\begin{aligned}
 f(x_i; \mu, \sigma) &= \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x_i - \mu)^2 / (2\sigma^2)} \\
 L(\mu; \mathbf{x}) &= \prod_{i=1}^n f(x_i; \mu, \sigma) \\
 l(\mu; \mathbf{x}) &= \log(L(\mu; \mathbf{x})) \\
 &= \sum_{i=1}^n \log(f(x_i; \mu, \sigma)) \\
 &= \sum_{i=1}^n \left[ \log\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right) + (-(x_i - \mu)^2 / (2\sigma^2)) \right]
 \end{aligned}$$

Vi finner verdien av  $\mu$  som maksimerer log-likelihooden ved å derivere log-likelihooden, sette dette uttrykket lik 0, og løse for  $\mu$ . Mer presist er maksimum-likelihood-estimatoren  $\hat{\mu}$  for  $\mu$  den verdien av  $\mu$  som gir

$l'(\mu, \mathbf{x}) = 0$ , det vil si at  $\hat{\mu}$  er løsningsen av  $l'(\hat{\mu}, \mathbf{x}) = 0$ :

$$l'(\mu; \mathbf{x}) = \sum_{i=1}^n [(-2(x_i - \mu)/(2\sigma^2))]$$

$$l'(\hat{\mu}, \mathbf{x}) = 0$$

$$\sum_{i=1}^n [(-2(x_i - \hat{\mu})/(2\sigma^2))] = 0$$

$$\sum_{i=1}^n (x_i - \hat{\mu}) = 0$$

$$\sum_{i=1}^n x_i - n\hat{\mu} = 0$$

$$\sum_{i=1}^n x_i = n\hat{\mu}$$

$$\frac{\sum_{i=1}^n x_i}{n} = \hat{\mu}$$

$$\bar{x} = \hat{\mu}$$