

2. obligatoriske oppgave i STK 3100 høsten 2009.

Utlevering: Fredag 23. oktober.

Innleveringsfrist: Fredag 6. november kl. 14.30.

Besvarelsen innleveres i ekspedisjonen i 7. etasje, Niels Henrik Abels hus.

Dette er det andre settet med obligatoriske innleveringer i STK 3100 høsten 2009. Oppgavesettet består av 1 oppgave. Det er valgfritt om du vil skrive besvarelsen for hånd eller om du vil bruke et tekstbehandlingsprogram. Der du bruker R (eller et annet program), må utskrifter legges ved eller limes inn. Hvis flere studenter samarbeider om å løse oppgavene, må likevel hver student levere sin selvstendige besvarelse. Det må gå fram av besvarelsen hvem du har samarbeidet med. Se ellers ”Regelverk for obligatoriske oppgaver” som er gitt på kursets hjemmeside.

Obligatorisk oppgave:

Oppgaven er en modifikasjon og utvidelse av Oppgave 14 gitt til torsdag 22. oktober.

I denne oppgaven skal vi analysere noen data fra Baxter et al. (1980) *Transactions 21 Congress of Actuaries* **2-3**, 11-29 om skadetilfeller i en portefølje av forsikrede privatbiler i et middels stort engelsk forsikringsselskap tredje kvartal 1973. Det er registrert antall skadetilfeller oppdelt etter tre tarifferingsfaktorer, hver med fire nivåer. Tarifferingsfaktorene er kodet som følger:

- Forsikringstagers alder:
 - 1 = under 25 år.
 - 2 = 25-29 år.
 - 3 = 30-35 år.
 - 4 = over 35 år.

- Bilens motorvolum:
 - 1 = under 1 liter.
 - 2 = 1-1,5 liter.
 - 3 = 1,5-2 liter.
 - 4 = over 2 liter.

- Distrikt:
4 = London og andre store byer.
1–3 = andre distrikt.

Dataene er lagt på filen

<http://www.math.uio.no/avdc/kurs/STK3100/Data/claims>

slik at første søyle i filen angir alder, andre søyle bilens motorvolum, tredje søyle distrikt, fjerde søyle antall forsikrede i gruppen og femte søyle antall skader.

- Anta at for hver forsikringstaker inntreffer skader med en rate (som kan avhenge av forsikringstakers alder og bosted/distrikt og bilens motorvolum). Hvorfor er det rimelig å tenke seg at antall skader angitt i datafilen er Poissonfordelt med forventning proporsjonalt med antall forsikringstakere.
- Vis at kanonisk link for Poissondata er $g(\mu) = \log(\mu)$ (der μ er forventningen i en Poissonfordeling). Vi skal videre i oppgaven bare bruke denne linkfunksjonen.
- For de foreliggende dataene er altså forventet antall skader proporsjonalt med antall forsikringstakere. Forklar hvorfor dette antallet (med en loglink) inngår i lineær prediktor som en 'offset'.
- Diskuter om de tre kovariatene bør modelleres som kategoriske forklaringsvariable (faktorer) eller numeriske (kvantitative) forklaringsvariable. Diskuter fordeler og ulemper ved begge tilnærminger.
- Sett opp en deviansanalysetabell for dataene der alder, motorvolum og distrikt er modellert som faktorer. Forklar hvorfor modellen med alle hovedeffekter, men ingen interaksjoner, er adekvat for dataene.
- Forklar hvorfor parametrene i modellen har fortolkning som logaritmen til rateratioer (RR). Estimer rateratioene og beregn tilhørende konfidensintervall.
- Påvis at det er en lineær trend i alder og motorvolum ved å tilpasse modeller der disse tarifferingsfaktorene tas med som kvantitative kovariater ('variates' i GLM terminologi). Vis at det er mulig å foreta en forenkling i modelleringen av effekten av distrikt. (Dvs. det modelleres med færre parametere.) Tilpass denne forenklete modellen og fortolk resultatene. Gjør også en test for forskjell med modellen der variablene inngikk som faktorer.
- Foreta en analyse av residualene i den 'endelige modellen' du har kommet fram til over. Er det noe ved residualene som tyder på at denne ikke er tilfredsstillende?

-
- i) Estimer raten for skadetilfeller i 3. kvartal for en forsikret i alder 25-29 år med bil med motorvolum 1,5-2 liter bosatt i London. Finn også 95% konfidensintervall for denne raten.
- j) Undersøk om det er overspredning (evt. underspredning) i forhold til Poissonfordeling i dette datasettet.
- k) Vi har så langt benyttet log-link der logaritmen til antall forsikringstakere inngår som en offset. For andre linkfunksjoner må en da gå litt annerledes til verks. Først vis at når $Y \sim \text{Po}(n\lambda)$ blir

$$E[Y/n] = \lambda \quad \text{og} \quad \text{Var}[Y/n] = \frac{\lambda}{n}.$$

Lag nå en avledet responsvariabel lik antall skader delt på antall forsikringstakere. Tilpass så en Poissonregresjonsmodell med denne responsvariabelen hvor du også legger inn antall forsikringstakere som vekter: `weights=antforsikrede`. Sammenlign estimator og devians med analysen der dette antallet ble lagt inn ved 'offset'.

- l) Denne metoden å tilpasse dataene (dvs. med den avledede variabelen vektet med andre forsikringstakere) kan også brukes for andre linkfunksjoner, f. eks. potenslinkene $g_\rho(\mu) = \mu^\rho$. Forklar hvorfor loglinken kan betraktes som et grensetilfelle når $\rho \rightarrow 0$. Finn tilnærmet MLE for ρ (bruk `family=poisson(link=power(rho))` for passende valg av `rho`). Test også om loglinken er adekvat for dataene.