

Notat 2:

Betinget logistisk regression.

Udgangspunktet er en logistisk regressionsmodel for binære data, som vi kan skrive

$$P(y_{gi} = 1) = p_{gi}, \quad g = 1, \dots, G, \quad i = 1, \dots, n_g$$

hvor

$$p_{gi} = \frac{\exp(\alpha_g + \beta x_{gi} + \dots)}{1 + \exp(\alpha_g + \beta x_{gi} + \dots)} \quad \text{eller} \quad \text{logit}(p_{gi}) = \alpha_g + \beta x_{gi} + \dots$$

Her symboliserer “ $\beta x_i + \dots$ ” en modelformel, der i princippet kan være hvad som helst der kan stå på højre side af lighedstegnet i (f.eks.) en lineær normalfordelingsmodel.

Vi har her valgt en dobbelt indicering af observationerne y_{gi} med et gruppenummer g og et løbende nummer i inden for gruppen, for at markere den særlige rolle som faktoren “gruppe” har i det følgende. Bemærk at modellen involverer en hovedvirkning af denne faktor. I praksis må man forestille sig, at det drejer sig om en faktor med mange niveauer, som kun er medtaget af nød, idet det man primært er interesseret i ikke er forskelle mellem grupper, men de øvrige parametre β, \dots .

For at slippe af med parametrene α_g kan man betinge med summerne y_g inden for grupper. Denne betingning vil ikke fjerne væsentlig information, idet gruppsummerne først og fremmest indeholder information om parametrene α_g .

SÆTNING. I den betingede fordeling af y_{gi} 'erne, givet summerne y_g , bliver likelihoodfunktionen

$$\prod_{g=1}^G \left(\frac{\exp\left(\sum_{i:y_i=1}(\beta x_{gi} + \dots)\right)}{\sum_M \exp\left(\sum_{i \in M}(\beta x_{gi} + \dots)\right)} \right)$$

hvor summen i nævneren skal tages over de delmængder M af $\{1, \dots, n_g\}$ der netop har y_g elementer.

Beviset er forholdsvis elementært, men meget besværligt at skrive op. Man ser, at når man danner den betingede sandsynlighed for at få en bestemt konfiguration af ettaller inden for en gruppe, givet at der netop

skal være y_g . af dem, så er der en forfærdelig masse ting der forkorter ud. Først og fremmest forkorter alle nævnerne i udtrykkene for de enkelte sandsynligheder og deres komplementære ud. Derefter forkorter parametrene α_g ud, fordi de forekommer netop y_g . gange i tælleren (under exp-tegnet) og også netop y_g . gange i hvert af leddene i nævneren. Beviset overlades i øvrigt til læseren.

I ISUW kan kommandoen FITCLOGIT håndtere disse modeller, se ISUW-eksemplet CLOGIT.ISU til opgave 2, MPAS efteråret 2001. Jeg har ikke kendskab til andre programmer som kan håndtere disse modeller. Hvis man bruger SAS må man nøjes med at fitte den ubetingede model.