

## 19. Logistische regressie

### Achtergrond

Logistische regressie is multiële regressie waarbij de uitkomstvariabele categorisch is, en de predictorvariabelen continu of categorisch zijn. Simpel gezegd voorspelt logistische regressie in welke categorie personen vallen op basis van andere informatie. Als de uitkomstvariabele twee categorieën heeft, wordt het *binair logistische regressie* genoemd en bij meerdere categorieën wordt het *multinomiale logistische regressie* genoemd.

### De achterliggende principes

Een gewone regressieanalyse kan niet gebruikt worden omdat bij een categorische uitkomstvariabele de assumptie dat er een lineaire relatie is tussen de variabelen geschonden wordt. Een niet lineaire relatie kan door een logaritmische transformatie (de logit) alsnog lineair worden gemaakt. Dat is wat je doet bij logistische regressie.

Bij logistische regressie voorspel je niet de waarde van de uitkomstvariabele Y vanuit de predictor X, maar bereken je de kans op Y bij de waarden van X. Wanneer je maar één predictorvariabele hebt, kun je de kans op Y uitrekenen met deze formule:

$$P(Y) = \frac{\text{Odds}_{\text{stans na woedsel}}}{\text{Odds}_{\text{stans na affectie}} + 1}$$

$P(Y)$  is de kans dat Y voorkomt (de kans dat iemand in een bepaalde categorie hoort) en staat voor de natuurlijke logaritmes. Wanneer er meerdere voorspellers zijn, wordt het gedeelte binnen de haakjes uitgebreid met een extra predictor en de bijbehorende parameters.

De uitkomst van deze formule varieert tussen de 0 en de 1, waarbij waarden dichtbij 0 betekenen dat Y heel onwaarschijnlijk is. Een waarde dichtbij 1 zegt juist dat Y heel waarschijnlijk is.

De waarden van de parameters worden geschat met de *maximum-likelihood schatting*, zoals er in lineaire regressie de least squares methode wordt gebruikt. Deze schatting selecteert de coëfficiënten die het meest overeenkomen met de geobserveerde gegevens.

Het model testen: de log-likelihood statistiek

$P(Y_i)$  staat voor de kans op Y bij de  $i^{\text{ste}}$  persoon. Voor een bepaald persoon is Y (de werkelijke uitkomst) 0 of 1, de uitkomst vond wel plaats of de uitkomst vond niet plaats. Als je wil kijken hoe goed de fit van een model is, kun je de geobserveerde uitkomsten vergelijken met de voorspelde uitkomsten. Bij logistische regressie bepaal je de fit met de *log-likelihood statistiek*.

$$\text{Log-likelihood} = \sum [Y_i \ln(P(Y_i)) + (1 - Y_i) \ln(1 - P(Y_i))]$$

De log-likelihood komt ongeveer overeen met  $SS_R$  in multiple regressie. Ze zijn allebei een indicator van hoeveel onverklaarde informatie er is nadat het model en de gegevens vergeleken zijn. Hoge waarden van de log-likelihood statistiek betekent een slechte fit.

## Het model testen: de deviantie statistiek

Deze statistiek lijkt erg op de log-likelihood statistiek, het is namelijk  $-2 \times \log$ -likelihood. De deviatie statistiek wordt daarom vaak de  $-2LL$  genoemd. Het grote voordeel is dat deze statistiek een chi-square distributie heeft.

De log-likelihood kan voor meerdere modellen uitgerekend worden en de verschillende modellen kunnen dan vergeleken worden door te kijken naar het verschil tussen de deviaties. Het is zinvol om een model te vergelijken met een baseline. Bij lineaire regressie gebruiken we het totale gemiddelde als basis waarmee we het model vergeleken. Bij een categorische uitkomst betekent een gemiddelde niets, dus dat is niet zo zinvol. Daarom wordt bij logistische regressie de uitkomst die het vaakst voorkomt als baseline gebruikt. Dit is het logistische model als alleen de constante als model wordt gebruikt. Je kunt kijken wat de verbetering is van het toevoegen van voorspellers door de modellen te vergelijken:

$$X^2 = (-2LL(\text{baseline})) - (-2LL(\text{nieuw}))$$

$$X^2 = 2LL(\text{nieuw}) - 2LL(\text{baseline})$$

$$df = k_{\text{nieuw}} - k_{\text{baseline}}$$

Dit verschil heet een likelihood ratio en het heeft een chi-square distributie, waarbij het aantal vrijheidsgraden het aantal parameters in het nieuwe model – het aantal parameters in het baseline model is. Het baseline model bevat alleen de constante, dus heeft slechts 1 parameter.

*R en  $R^2$*

Het logistische equivalent van R bij lineaire regressie is de R-statistiek. De R-statistiek is de partiële correlatie tussen de uitkomstvariabele en elk van de voorspellervariabelen en het kan variëren tussen de -1 en 1. Een positieve waarde betekent dat wanneer de predictorvariabele toeneemt, de kans op de gebeurtenis ook toeneemt. Een negatieve R betekent dat de kans afneemt als de predictorvariabele toeneemt. Wanneer een variabele een lage R heeft draagt het weinig bij aan het model.

$$R = \sqrt{\frac{z^2 - 2df}{-2LL(\text{baseline})}}$$

De  $-2LL$  is de log-likelihood voor het originele model, de  $z$  staat voor de Wald statistiek en de  $df$  kun je aflezen in de samenvattende tabel bij de variabelen in de vergelijking. R moet met voorzichtigheid gebruikt worden, omdat de Wald statistiek niet altijd accuraat is en deze waarde mag niet gekwadeerd worden, zoals bij  $R^2$  in lineaire regressie.

Er zijn meerdere statistieken die ongeveer equivalent zijn aan  $R^2$ , zoals *Hosmer and Lemeshow's  $R^2_L$* , Cox en Shell's  $R^2_{CS}$  (deze wordt door SPSS gebruikt, maar haalt nooit het theoretische maximum) en Nagelkerke's  $R^2_N$ . De formules zijn te vinden op pagina 765 en 766 van het boek. Hoewel ze verschillende antwoorden opleveren, zijn ze allemaal ongeveer gelijk aan het idee van  $R^2$ .

*De contributie van voorspellers bepalen*

Om de contributie van de voorspellers aan het model te bepalen, wordt in logistische regressie de z-statistiek gebruikt, welke een normale verdeling heeft. Het is vrijwel gelijk aan de t-toets die in lineaire regressie gebruikt wordt om de contributie van voorspellers te meten.

Deze statistiek vertelt ons of de  $b$  coëfficiënt voor de voorspeller significant verschilt van 0, en dus een significante bijdrage levert aan het model.

$$z = \frac{b}{SE_b}$$

De  $z$ -statistiek wordt de *Wald statistiek* genoemd. SPSS geeft  $z^2$ , welke een chi-square verdeling heeft. Je moet echter uitkijken met het gebruik van de Wald statistiek, want bij een grote  $b$  is er een inflatie van de standaard meetfout ( $SE_b$ ), waardoor de  $z$  te laag uitvalt en er een risico is op een type II fout. Je kunt vaak beter om de voorspellers op de hiërarchische manier in het model te plaatsen en te kijken naar de verandering in likelihood ratio statistieken om te bepalen of een voorspeller een goede bijdrage levert.

#### *De odds ratio*

Voor het interpreteren van logistische regressie is de odds ratio het belangrijkste. Dit is de exponent van  $b$  ( $e^b$  of  $\exp(B)$ ). Het geeft de verandering in kans dat komt door de verandering in de predictor. De kans dat een gebeurtenis vóórkomt, gedeeld door de kans dat een gebeurtenis niet vóórkomt, is de *odds*. Als voorbeeld wordt genomen dat we voorspellen of iemand zwanger is geraakt door te kijken naar of er een condoom werd gebruikt tijdens de laatste geslachtsgemeenschap. De odds om zwanger te worden is de kans om zwanger te worden gedeeld door de kans om niet zwanger te worden.

$$\text{Odds} = \frac{P(\text{gebeurtenis})}{P(\text{geen gebeurtenis})}$$

$$P(\text{gebeurtenis } Y) = \frac{1}{1 + e^{-(b_0 + b_1 X_1)}}$$

$$P(\text{geen gebeurtenis } Y) = 1 - P(\text{gebeurtenis } Y)$$

Om de verandering in odds te berekenen wanneer de predictor verandert, moet je de odds berekenen om zwanger te raken wanneer je wel een condoom gebruikt en wanneer je geen condoom gebruikt. De odds ratio is de proportionele verandering in odds na een verandering in de predictor:

$$\text{Odds ratio} = \frac{\text{odds na verandering van een unit in de voorspeller}}{\text{originele odds}}$$

Wanneer de uitkomst een grotere waarde is dan 1, betekent het dat als de voorspeller toeneemt, de odds van de uitkomst ook toeneemt.

#### *Het maken van een model en parsimonie*

Als je meer dan één voorspeller hebt, kun je kiezen uit dezelfde methodes om een model te maken als bij lineaire regressie. Net zoals bij gewone regressie zijn de forced entry en hiërarchische methode aan te bevelen. De stepwise methode is ook hier af te raden, maar als je het echt per se wil, je hebt geen theorieën en causaliteit je niet uit maakt, kun je beter de backward methode gebruiken. Dit is vanwege de *onderdrukkende effecten* die de forward methode kan hebben. Deze komen voor wanneer een voorspeller een significant effect heeft maar alleen wanneer een andere variabele constant wordt gehouden. Ook kun je beter de likelihood ratio methode kiezen en niet de Wald statistiek, omdat die in bepaalde situaties erg onbetrouwbaar is.

Het is net zoals bij lineaire regressie het beste om de hiërarchische methode te gebruiken en een model te maken op een systematische manier, op basis van theorieën. Daarnaast is er een streven naar *parsimonie*, wat betekent dat een eenvoudige verklaring van een fenomeen te verkiezen is boven een ingewikkelde verklaring. Dat betekent dat er dus geen predictors in een model moeten zitten, tenzij ze van belang zijn in het verklaren van het fenomeen.

Belangrijk om hierbij te onthouden is wel dat interactietermen alleen geldig zijn als de hoofdeffecten ook in het model zitten. Ook als de hoofdeffecten zelf niet veel toevoegen mogen ze niet verwijderd worden als de interactieterm in het model zit.

## Assumpties en dingen die fout kunnen gaan

De assumpties voor logistische regressie zijn grotendeels hetzelfde als besproken in hoofdstuk 5 en 8. Bijzonder om op te merken zijn de volgende twee assumpties:

**Lineariteit:** Hier betekent het dat er een lineaire relatie moet zijn tussen de continue voorspellers en de logit van de uitkomstvariabele. Je test deze assumptie door te kijken of het interactie-effect tussen de predictor en de logtransformatie significant is.

**Onafhankelijkheid van de meetfouten:** Schenden van deze assumptie zorgt bij logistische regressie voor overspreiding. Hierover wordt later meer verteld.

Logistische regressie heeft ook specifieke problemen en SPSS kan door die problemen een totaal verkeerde output geven. SPSS werkt met ‘iterations’, het schat de parameters. Met elke nieuwe poging (iteration) wordt geprobeerd een accuratere schatting te maken. Het stopt met schatten als het maximale aantal schattingen dat je hebt ingevoerd is bereikt, of wanneer convergentie is bereikt. Dit laatste betekent dat elke nieuwe schatting vrijwel hetzelfde antwoord oplevert. In twee situaties kan het voorkomen dat convergentie niet wordt bereikt, wat een volledig incorrecte output oplevert: incomplete informatie en complete scheiding.

**Incomplete informatie van de predictors** betekent dat je niet van alle combinaties van de variabelen data hebt verzameld. SPSS kan in dat geval geen goede uitkomst geven. Je kunt dit controleren door een contingency tabel te maken via Crosstabs, zoals beschreven in het vorige hoofdstuk. Je kunt het ook zien aan coëfficiënten die een buitengewoon grote standaarderror hebben. Je kunt het oplossen door meer data te verzamelen.

Een tweede probleem kan ontstaan als de uitkomstvariabele perfect wordt voorspeld door een voorspeller of een combinatie van voorspellers. Dit heet *complete scheiding*. Er zijn dan bijvoorbeeld twee categorieën waarbij de data niet overlappen. Een voorbeeld is het voorspellen of iets een kat of een inbreker is op basis van gewicht. Alles onder de 15 kilo is een kat, alles boven de 40 kilo is een inbreker, er is geen overlap in gewicht tussen katten en inbrekers. SPSS weet dan niet hoe het de tussenruimte waar geen data is (tussen 15 en 40 kilo) moet invullen. Dit probleem ontstaat vaak wanneer er te veel variabelen worden gebruikt voor te weinig data. Een oplossing is dan meer data verzamelen. Soms kan een simpeler model uitkomst bieden.

### *Overspreiding (overdispersion)*

Er is sprake van overspreiding wanneer de geobserveerde variantie groter is dan de verwachte variantie uit een logistisch regressiemodel. Dit kan veroorzaakt worden doordat de aanname van onafhankelijkheid niet wordt geschonden of doordat de kans op succes varieert. De standaardmeetfouten zijn dan te klein, waardoor de teststatistiek te groot is en er te snel statistische resultaten worden gevonden.

Ook het betrouwbaarheidsinterval is dan te smal, waardoor we verkeerde aannames kunnen doen over de populatie.

Overspreiding is aanwezig wanneer de ratio van de chi-square goodness-of-fit statistiek en het aantal vrijheidsgraden groter dan 1 is (deze ratio heet de dispersion parameter,  $\Phi$ ). Het wordt problematisch als deze ratio nabij of groter dan 2 is.

De effecten van overspreiding kunnen worden verminderd door de dispersion parameter te gebruiken om de standaardmeetfouten en de betrouwbaarheidsintervallen aan te passen. De standaarderror kan

bijvoorbeeld vermenigvuldigd worden met  $\sqrt{\Phi}$ .

## Binaire logistische regressie met een waargelijk voorbeeld

Een man werd opgenomen in het ziekenhuis met een geperforeerd rectum, omdat hij een levende aal ingebracht had in zijn anus. De verklaring van de man was dat hij dacht dat het zou helpen tegen obstipatie. Om deze hypothese te testen, kun je data verzamelen, met als uitkomstvariabele of de obstipatie is genezen of niet. Er is dus een dichotome uitkomstvariabele. De eerste predictorvariabele is de interventie, de aal in de anus of geen behandeling. De tweede predictorvariabele is de duur, het aantal dagen dat de obstipatie al duurt.

Er zijn drie mogelijke modellen. Het eerste model heeft alleen de interventie als predictor, de tweede heeft ook de duur als predictorvariabele en in het derde model is ook de interactie interventie x duur toegevoegd. Er moet gekeken worden welke van de drie modellen de beste fit heeft, terwijl je ook rekening houdt met parsimonie.

Als besloten is welke van de drie modellen het beste is, doe je de analyse opnieuw voor dit specifieke model, met de diagnostische statistieken erbij wie je kunt inspecteren op bronnen van bias, zoals uitschieters en invloedrijke gevallen. Daarna controleer je de lineariteit van de logit en multicollineariteit.

De categorische variabelen worden in SPSS gecodeerd met 1 (wel gebeurtenis) en 0 (geen gebeurtenis). In het voorbeeld krijgen de categorieën wel genezen en wel interventie een 1, niet genezen en geen interventie krijgen de code 0.

In SPSS ga je voor het maken van het model naar Analyze - Regression - Binary logistic. Bij Dependent zet je de uitkomstvariabele. Je wil de drie modellen invoeren als drie blokken. Om een hoofdeffect van een variabele in het model op te nemen, selecteer je een predictor en sleep je die naar Covariates. Voor het invoeren van een interactie-effect selecteer je meerdere predictorvariabelen tegelijk en voer je ze in bij Covariates. Door op de knop Next te klikken, krijg je een nieuw blok waar je het volgende model kunt invoeren.

Bij Categorical kan je aangeven welke voorspellers categorisch zijn. SPSS gebruikt standaard Indicator codering, wat gewone dummy codering is, waarbij je kunt kiezen om de eerste of de laatste categorie als baseline te gebruiken.

Je kunt nu de analyse uitvoeren zonder de andere opties te gebruiken, om te kijken welk model de beste fit heeft. In de output verschijnen nu de samenvattingen van de modelstatistieken voor de drie modellen. In de rij Model vind je de chi-square statistiek en significantie voor het hele model, bij Block zie je de verandering sinds het vorige model. Als het eerste model significant is, ben je daarna alleen geïnteresseerd in de significantie bij Block, die zegt of het nieuwe model een verbetering is boven de vorige. In het aalvoorbeeld wordt model 1 gekozen, waarbij alleen interventie een predictor is.

Daarna kun je de analyse opnieuw doen, met alleen het eerste model. Bij de optie save kan je SPSS variabelen met residuen aan laten maken om te kijken hoe goed het model bij de gegevens past. De meeste opties zijn bekend vanuit de gewone regressie, alleen bij logistische regressie heb je ook nog de Predicted probabilities en de Predicted group memberships. Predicted probabilities is de kans dat Y voorkomt bij de waarden van elke voorspeller voor een bepaalde deelnemer. Predicted group memberships voorspelt tot welke van de categorieën de deelnemer hoort, gebaseerd op het model.

Field adviseert om in elk geval Probabilities, Group membership, Cook's, Leverage values, DFBeta(s), Standardized en Include the covariance matrix aan te klikken. De variabelen worden niet opgeslagen als je bootstrap gebruikt.

Bij Opties zijn de standaard instellingen meestal prima. Classification plots geeft histogrammen van de werkelijke en voorspelde waarden van de uitkomstvariabelen, wat handig kan zijn om de fit van het model te bekijken. Casewise listing of residuals geeft de gevallen waarbij de standardized residual meer dan 2 standaardafwijkingen afwijkt. De Hosmer-Lemeshow goodness-of-fit statistiek wordt gebruikt om de fit van het model te testen. Iteration history is handig om aan te klikken omdat je daarmee de eerste waarde van  $-2LL$  krijgt, die je nodig hebt om R uit te rekenen.

Bootstrap kun je alleen gebruiken als je de forced entry methode gebruikt. Het slaat de variabelen bij save niet op, dus als je het beide wil, moet je de analyse opnieuw uitvoeren.

## Interpreteren van logistische regressie

### *Blok 0*

De output is in twee blokken verdeeld. Blok 0 geeft het model voordat de predictor is toegevoegd, blok 1 geeft het model met de predictor interventie erbij. Het enige dat nuttig is in blok 0 is de Iteration history, waar je de initiële  $-2LL$  kunt aflezen.

### Het model

Bij Model summary staat de  $-2LL$  voor het model. In de Classification table zie je hoe goed het model groepslidmaatschap voorspelt. Het model gebruikt interventie als predictor, waarbij het mensen die de interventie ondergingen classificeert als genezen en de mensen zonder interventie classificeert als niet genezen.

De tabel Variables in the equation is erg belangrijk, het geeft de coëfficiënten en statistieken voor de predictors in het model. In de SPSS output is de Wald statistiek gekwadrateerd, en het test of de b-coëfficiënt significant afwijkt van 0. Voor de data in het voorbeeld is de Wald statistiek significant, wat betekent dat interventie een significante predictor is voor of de patiënt genezen is.

De R-statistiek en  $R_L^2$  zijn met de hand te berekenen aan de hand van de output. Cox en Snell's en Nagelkerke's metingen staan in de SPSS output.

De odds ratio staat in de SPSS output als  $\text{Exp}(B)$ . Waardes groter dan 1 betekenen dat bij een toename van de predictor de odds van de uitkomstvariabele toenemen. Waardes kleiner dan een betekenen een afname van de odds van de uitkomstvariabele bij toename van de predictor. Bij het betrouwbaarheidsinterval kijk je dus of de 1 in het interval zit. Als dit zo is, kan het effect beide kanten op gaan. In het voorbeeld betekent dat dat interventie kan leiden tot meer kans op genezing of tot minder kans op genezing.

De classificatieplot laat een grafiek zien met de voorspelde kansen dat een patiënt genezen is. Als het model een perfecte fit heeft, zie je alle genezen patiënten rechts en alle niet genezen patiënten links.

Dus, hoe meer gevallen er op de uiteinden van de grafiek zitten, hoe beter het is. Een kans van 0.5 betekent dat het model even goed is als puur toeval. Je hoopt dus dat de gevallen ver van die 0.5 af zitten.

## Voorspelde kansen

SPSS kan een lijst geven van de voorspelde kansen van de uitkomstvariabele op basis van het model. Met de optie Save heb je de variabelen voorspelde kansen en voorspelde groepslidmaatschappen gemaakt. Die variabelen kun je als lijsten in de output krijgen via Analyze - Reports - Case summaries.

### *Interpreteren van de residuen.*

Je controleert de residuen in een regressieanalyse om te kijken of er punten zijn waarop het model een slechte fit heeft en om te kijken naar gevallen die erg veel invloed uitoefenen op het model. Om de fit te controleren gebruik je de residuen, specifiek de Studentized residual, standardized residual en deviatie-statistieken. Voor het tweede punt gebruik je invloedsstatistieken zoals Cook's distance, DFBeta en leverage statistieken. Deze zijn in hoofdstuk acht uitgebreid besproken, een samenvattende tabel staat op pagina 791.

Voor de effectgrootte kan de odds ratio het best gebruikt worden.

## Rapporteren van logistische regressie

Logistische regressie kan je ongeveer hetzelfde rapporteren als lineaire regressie. Vermeld de b-waardes, standaard meetfouten en de significantie. Vermeld ook de odds ratio, het betrouwbaarheidsinterval en de constante.

## Testen van assumpties

Voor het testen van de assumptie van lineariteit voer je de logistische regressieanalyse opnieuw uit, maar dan voeg je predictors toe die de interactie zijn van elke predictor en de log van zichzelf. Van de log van de variabele moet je een nieuwe variabele maken. Dit doe je via Transform – Compute Variable.

Om de assumptie te testen voer je de analyse hetzelfde uit als eerst, maar dan voeg je alle variabelen (ook de nieuwe interactietermen) in één keer toe. Je gebruikt dus forced entry in plaats van een hiërarchische methode. Significante interactietermen betekenen een schending van de assumptie.

Bij logistische regressie is multicollineariteit ook een probleem, maar via de logistische analyse kan de assumptie niet getest worden. Het testen moet daarom via een lineaire regressieanalyse, met dezelfde uitkomstvariabele en predictorvariabelen (dus ook de interactietermen). Vink bij Statistics de Collinearity diagnostics aan voor multicollineariteit.

Tolerantiewaarden kleiner dan 0.1 en VIF waarden groter dan 10 zijn problematisch. In de tabel Collinearity Diagnostics staan eigenwaardes. Wanneer de eigenwaardes ongeveer even groot zijn, betekent het dat kleine veranderingen in de variabele het model niet veranderen. Als een van de eigenwaardes veel groter is dan de rest, heb je een probleem. Als de condition indexes veel groter zijn dan 1, geeft dat ook multicollineariteit aan. Tenslotte kun je de variantie proporties bekijken. Als een voorspeller met een kleine eigenwaarde een grote proportie variantie heeft, is er sprake van collineariteit.

Een echte oplossing voor multicollineariteit is er niet. Je kan de variabele uit het model halen of de steekproef vergroten. Maar het beste is erkennen dat het model niet betrouwbaar is.

## Voorspellen van meerdere categorieën: multinomiale logistische regressie

*Multinomiale logistische regressie* is het gebruiken van logistische regressie om groepslidmaatschap van meer dan twee categorieën te voorspellen. Hierbij wordt de uitkomstvariabele gesplitst in een aantal vergelijkingen tussen twee categorieën. Je moet hierbij een baseline categorie kiezen waarmee je alles vergelijkt.

### *Multinomiale logistische regressie in SPSS*

Voor deze vorm van regressie in SPSS ga je naar Analyze - Regression - Multinomial logistic. De uitkomstvariabele komt bij Dependent, de categorische voorspellers bij Factors en de continue voorspellers bij Covariates. Bij Reference category kan je aangeven of je alle categorieën wil vergelijken met de eerste of de laatste categorie.

Interactie-effecten specificeren kan bij multinomiale logistische regressie bij Model, waarbij je kiest voor Custom model. De variabelen kunnen in forced entry terms of stepwise terms. Zonder de hoofdeffecten zou een deel van de variantie niet verklaard worden dus deze worden ook toegevoegd, met forced entry. De interacties kunnen eventueel wel bij stepwise ingevoerd worden, dan komen de interacties alleen in het model als ze significante voorspellers zijn.

### *Statistieken*

Bij Statistics komen er een aantal statistieken naar voren:

De pseudo R-square geeft Cox-Snell en Nagelkerke's  $R^2$ . Deze kunnen als effectgrootte worden gebruikt. Step summary vat elke stap samen met verwijderde en toegevoegde voorspellers. Het wordt gebruikt als je kiest voor de stepwise methode. Model fitting information vergelijkt het model met de basis, een model zonder predictors, met alleen de intercept. Information criteria is handig voor het vergelijken van modellen als je de stepwise methode gebruikt of verschillende modellen wil vergelijken. Cell probabilities geeft een tabel van de geobserveerde en verwachte frequenties. Classification table laat de geobserveerde versus de verwachte waarden zien voor alle combinaties van de voorspellers. Deze optie kun je beter niet gebruiken bij grote analyses. Goodness-of-fit geeft de Pearson en de likelihood ratio chi-square voor het model en is dus belangrijk om te selecteren. Monotonicity measures is alleen handig als je uitkomstvariabele slechts twee uitkomsten heeft. Estimates geeft de b-waarden, test statistieken en het betrouwbaarheidsinterval voor de voorspellers in het model. Dit is een belangrijke optie. Likelihood ratio tests geeft dezelfde informatie als de significantiewaarden van de individuele voorspellers. Asymptotic correlations and covariances geeft de correlaties tussen de beta's in het model.

### *Andere opties*

Bij Criteria staan opties die je kunt gebruiken om de 'iterations' waarmee logistische regressie werkt, aan te passen. Deze opties moet je met rust laten, tenzij je de melding krijgt dat het convergeren mislukt is. Als je de opties aanpast, moet je erop bedacht zijn dat het model dat hieruit komt niet altijd betrouwbaar is.

Bij Options staat de optie Scale, waarmee de standaard meetfouten worden gecorrigeerd, wat nodig is als je overspreiding hebt.

### *Interpreteren van de output*

Bij logistische regressie kan er soms een waarschuwing tevoorschijn komen in de output. Deze waarschuwing betekent dat je niet voor alle combinaties van variabelen gegevens hebt. Dit is onvermijdelijk als je meerdere continue variabelen hebt en de waarschuwing kan genegeerd worden. Het is wel handig om naar de coëfficiënten te kijken of ze niet te grote standaard meetfouten hebben.



Daarna komt als eerst in de output de samenvatting van de stappen van de analyse naar voren, omdat gekozen is voor een stepwise analyse. Hier worden de hoofdeffecten en de significante interactie-effecten weergegeven. Daaronder worden de statistieken van het uiteindelijke model gegeven. Deze tabel laat ook de likelihood ratio test zien over het hele model.

Het volgende deel van de output geeft de goodness-of-fit weer. Wanneer de Pearson en deviance statistiek niet significant zijn, dan zijn voorspelde waardes niet significant verschillend van de geobserveerde waardes en heeft het model dus een goede fit. In het voorbeeld geeft de Pearson een significant resultaat, maar de deviance niet. Dit kan komen door overspreiding of het kan zijn dat de Pearson statistiek vertekend is doordat er zoveel lege cellen zijn, wat de waarschuwing vertelde.

De tabel Likelihood Ratio Tests wordt gebruikt om te kijken naar de significantie van de voorspellers in het model. Dit zegt alleen of de variabele een significante voorspeller is, maar niet wat het precieze effect is. Daarvoor moet je kijken naar de individuele parameter schattingen. Die tabel is in twee delen gesplitst omdat er steeds twee categorieën van de uitkomstvariabele worden vergeleken.

Het rapporteren van de resultaten is hetzelfde als bij een binaire logistische regressie, behalve dat de resultaten in een tabel gesplitst worden.