

## INTERACTIE-MODELLEN / MODERATIE-ANALYSE

Harry Ganzeboom

20 november 2014

19 februari 2015

Van interactie-effecten (of ‘moderatie’) spreken wanneer het effect van een X-variabele op Y afhangt van de waarden van een andere X-variabele, bijvoorbeeld de invloed van opleiding op inkomen verschilt tussen jongere (onervaren) en oudere (ervaren) werknemers. Het woord ‘interactie’ is niet zo vanzelfsprekend en in de gedragswetenschappen ook niet gebruikelijk: daar spreekt men bij voorkeur over een ‘moderator’-variabele – overigens ook geen vanzelfsprekende term. Bij sociologen en economen is “interactie” ingeburgerd.

Interactie is het gemakkelijkst op te sporen door de data in (twee) groepen op te delen en de regressiemodellen met elkaar te vergelijken. In het bovenstaande voorbeeld betekent dit dat je een groep van jongere en oudere werknemers onderscheidt en afzonderlijk het effect van opleiding op inkomen berekent. Hoewel deze procedure niet algemeen is aan te bevelen, kan het een goed begin van een analyse zijn. Heel handig in SPSS is dat je over het SPLIT FILE commando beschikt:

```
Recode age (lo thru 40=0) (40 thru hi=1) into dage.  
Sort cases by dage.  
Split file by dage.  
Regress /dep=income /enter=educ.  
Split file off.
```

Dit bespaart met name veel typewerk wanneer je de data in meer dan twee groepen wilt verdelen. De strategie heeft echter een paar nadelen:

- Het wordt gauw onoverzichtelijk en daarom kies je er al snel voor niet meer dan twee groepen te onderscheiden. Als de moderatorvariabele in feite veelwaardig is (zoals leeftijd) verlies je informatie.
- Je krijgt op deze manier geen informatie over een zeer dringende vraag: is het verschil tussen de twee groepen statistisch significant?

Om die redenen is het te verkiezen interacties te modelleren via multiplicatieve interactiemodellen.

Een multiplicatief interactiemodel voor het bovenstaande ziet er als volgt uit:

```
Compute educ_dage=educ*dage.  
Regress /dep=income /enter=dage educ educ_dage.
```

Je kunt de coëfficiënten van het regressiemodel op de volgende manier interpreteren:

- B0 Inkomen als zowel *dage* als *educ* gelijk aan 0 zijn, dus zeer laag opgeleide jongeren.
- B1 Het verschil in inkomen tussen jongeren (*dage*=0) en ouderen (*dage*=1) als *educ*=0.
- B2 Het verschil in inkomen tussen personen met +1 niveau opleiding voor jongeren (*dage*=0).
- B3 Hoeveel meer inkomen ouderen krijgen voor +1 niveau opleiding;

In feite impliceert het model precies dezelfde regressiemodellen die we bij de data-splitsing verkregen. De SE en T-waarde van de B3 coëfficiënt zijn extra informatie en geven direct antwoord op de vraag of het verschil in effect tussen de groepen statistisch significant is.

Er is echter in deze procedure eigenlijk geen reden *age* te dichotomiseren tot *dage*. Heel goed werkt ook:

```
Compute educ_age=educ*age.  
Regress /dep=income /enter=age educ educ_age.
```

Nu is de interpretatie van de regressiecoëfficiënten:

- B0 Inkomen als zowel *age* als *educ* 0 zijn, dus zeer laag opgeleide pasgeborenen.
- B1 Het verschil in inkomen tussen pasgeboren (*age*=0) en 1-jarigen (*age*=1) als *educ*=0.
- B2 Het verschil in inkomen tussen personen met +1 niveau opleiding voor 0-jarigen (*age*=0).
- B3 Hoeveel meer inkomen men krijgt voor +1 niveau opleiding als men 1 jaar ouder wordt.

Van deze coëfficiënten heeft eigenlijk alleen B3 een fatsoenlijke interpretatie, B0, B1 en B2 zijn absurd. Statistisch gezien zijn ze extrapolaties (nul-jarigen zitten vaak niet in de data en zijn in ieder geval geen werknemers) en hun SE en T-waarden toetsen nulhypothese die ons eigenlijk niet interesseren.

Om deze absurditeiten te voorkomen en wat begrijpelijker getallen te verkrijgen is het aan te bevelen in interactiemodellen variabelen zo te herschalen dat zij een meeteenheid verkrijgen waarin 0 een begrijpelijke waarde is en 1 een begrijpelijke eenheid is. Dat kan op verschillende manieren:

1. Dichotomiseren en de twee groepen met 0,1 coderen (zie boven).

2. De variabelen lineair transformeren tot een 0..1 bereik: trek de onderste waarde af en deel door het bereik. Bv. `compute agex=(age-18)/50.`
3. Door de variabelen in percentiel (proportiel) scores uit te drukken: `rank age / prop.`
4. Door de variabelen te centreren rondom hun gemiddelde: `compute agec=age-42.`
5. Door de variabelen in z-scores uit te drukken: `desc age / save.`

Strategie 2, 4 en 5 zijn lineaire transformaties: ze veranderen niet het intrinsieke model, maar alleen de meeteenheid. T-waarden van de interactieterm blijven ongewijzigd.

Strategie 1 en 3 zijn niet-lineaire transformaties en hebben ook consequenties voor de T-waarden. Mijn eigen voorkeur is strategie 3, omdat je daarmee ook uitbijters weghaalt en een zeer begrijpelijke meeteenheid overhoudt. De andere vier strategieën treft je vaker in de literatuur aan.

Strategie 1, 2 en 3 transformeren naar een 0-1 bereik en de uitkomsten kunnen dus zeer soortgelijk geformuleerd worden. Je doet dat door aan te geven wat het model betekent voor de 0-groep en hoe dat verschilt / wat het wordt in de 1-groep.

Strategie 5 transformeert naar een Z-bereik, met gemiddelde 0 en 1 standaarddeviatie als eenheid. De intercept B0 is dan het inkomen voor laagst opgeleiden ( $educ=0$ ) voor mensen met gemiddelde leeftijd, terwijl B1 staat voor de invloed van opleiding op inkomen voor mensen met een gemiddelde leeftijd. B3 staat nu voor hoe het effect van opleiding verandert als je een standaarddeviatie ouder wordt. Het bijzondere voordeel van deze formulering dat de hoofdeffecten niet of weinig veranderen door de introductie van het de interactieterm.

```
Regress /dep=income /enter=educ zage.
Regress /dep=income /enter=educ zage educ_zage.
```

Deze twee modellen leveren dezelfde B0, B1 en B2 op.

Het systeem laat zich eenvoudig uitbreiden naar een discrete moderator variabele met meer dan twee categorieën:

```
Recode age (lo thru 30=1) (30 thru hi=0) into age1830.
Recode age (30 thru 45=1) (18 thru 30=0) (45 thru hi=0) into age3045.
Recode age (lo thru 45=0) (45 thru hi=1) into age4564.
Compute educ_age1830=educ*age1830.
Compute educ_age3045=educ*age3045.
Compute educ_age4564=educ*age4564.
Regress /dep=income /enter=age1830 age3045 educ educ_age1845
educ_age3045.
```

## Verdere opmerkingen

1. Hoewel interactiemodellen drie regressiecoëfficiënten schatten is het belangrijk je te realiseren dat het nog steeds over twee variabelen gaat. Het is onjuist om over de interactieterm als een derde variabele te spreken. Als  $X_1$  en  $X_2$  constant worden gehouden, varieert  $X_1 * X_2$  niet.
2. Interactie-effecten zijn symmetrisch tussen de twee X-variabelen. B3 in het bovenstaande voorbeeld staat er zowel voor hoeveel een eenheid leeftijd het effect van opleiding verandert, maar ook hoe een eenheid opleiding het effect van leeftijd verandert.
3. Hoewel symmetrisch van aard, is het soms toch handig om een asymmetrische interpretatie te kiezen en de variabelen zo te benoemen dat ze bv. staan voor hoe leeftijdsverschillen het effect van opleiding veranderen.
4. Het is voor het interpreteren van interactiemodellen van groot belang dat je het 0-punt en de 1-heid van een variabelen goed kent – en ook voor alle lezers. Dit klemmt des te meer als je variabelen transformeert om de interpretatie ‘begrijpelijker’ te maken. Neem dus altijd een tabel op met descriptives (mean, std, min, max). Het is heel jammer dat SPSS die tabel niet levert bij regressie.
5. Veel data-analysten kijken graag naar interactiemodellen door ze stapsgewijs op te bouwen – dat helpt inderdaad bij de interpretatie. Je kunt de statistische toetsing dan ook ontlenen aan een F-test voor de toename van verklaarde variantie. Inhoudelijk levert dat dezelfde resultaten op. De specificatie van het bovenstaande kan zijn: `regr /dep=income /enter=dage educ /enter=educ_dage / des=def change.`
6. Af en toe hoor je iemand wel eens klagen over “multicollineariteit” als probleem bij interactie-modellen. Als je in SPSS collineariteitsstatistieken TOL en VIF laat berekenen, zijn deze vaak heel laag (TOL) en hoog (VIF). Het is echter niet iets om je zorgen over te maken. Deze statistieken berusten op de multiple correlatie tussen de X-variabelen en die is in dit geval zeer hoog.  $X_1 * X_1$  wordt vanzelfsprekend helemaal bepaald door  $X_1$  en  $X_2$  en het is alleen aan de lineariteitsbeperking van het correlatiemodel te danken dat  $R^2$  niet 1.0 is. Maar  $X_1 * X_2$  is niet een aparte variabele, het zijn precies dezelfde variabele als  $X_1$  en  $X_2$ .

Een behulpzaam tegengif kan zijn variabelen te Z-standaardiseren voor de analyse. Maar als je goed kijkt, zie je dat de statistische toetsing van interactietermen helemaal niet afhankelijk is van gekozen standaardisatiestrategie. De hoofdeffecten (B1 en B2) zijn dat wel – en met reden. Het is een goede oefening te bedenken waarom dat zo is.

7. Interactie-effecten behoeven niet altijd handmatig (via compute) te worden aangemaakt, er zijn ook spss-procedures die ze automatisch aanmaken: UNIANOVA als vervanging van REGRESS en LOGISTIC, NOMREG en PLUM voor logistische regressie. Bij deze programma's wil het echter nog wel eens een zoekplaatje worden bij discrete interacties.