

EEN TOEPASSING VAN EEN MODEL VOOR OVER- EN/OF
ONDERDISPERSIE BY BINOMIALE VARIATIE

Bas Engel

0 SAMENVATTING

Dit artikel geeft een toepassing van een model voor over- en/of onderdispersie ten opzichte van de variantie van een Binomiale verdeling. Parameterschattingen volgen uit een combinatie van maximum quasi-likelihood (McCullagh en Nelder, 1983) voor de locatie-parameters en een momentenmethode voor de over-/onderdispersie parameter (Williams, 1982). Voor de numerieke verwerking is het statistische pakket Genstat (Genstat 5, 1987) gebruikt.

Groep Landbouwwiskunde
Postbus 100
6700 AC Wageningen
tel.: 08370-19100

Met dank aan Ir. Joop te Brake (voor het beschikbaar stellen van de gegevens) en Willem Buist (voor assistentie bij de numerieke verwerking), beide medewerkers van het IVO 'Schoonoord' te Zeist.

Voor een onderzoek naar de effecten van actieve immunisatie tegen Androsteendion op de vruchtbaarheid van Texelaar ooiën, uitgevoerd op het Instituut voor Veeteeltkundig Onderzoek "Schoonoord" te Zeist, zijn 125 ooiën verdeeld in een controlegroep en een proefgroep. De dieren in de proefgroep zijn behandeld met het preparaat Fecundin. Van deze behandeling wordt na dekking een stijging van het aantal ovulaties en van het aantal embryos verwacht.

De dieren zijn ingedeeld in 4 leeftijdsklassen: $\leq 1/2$ jaar, $1/2-1 1/2$ jaar, $1 1/2-2 1/2$ jaar en $> 2 1/2$ jaar. Verder is er sprake van 2 dekperiodes. Er zijn dus 3 proeffactoren: behandeling, leeftijd en dekperiode met respectievelijk 2, 4 en 2 niveaus.

In dit artikel kijken we naar de analyse van het aantal embryos, conditioneel op het aantal ovulaties. De waarnemingen zijn gegeven in tabel A1 in de appendix. Vier dieren, die niet drachtig werden (waarvoor het aantal ovulaties dus gelijk aan nul is) en twee dieren, waarvan het aantal ovulaties niet bekend was, zijn niet in de analyse opgenomen.

2 HET MODEL

Stel voor een willekeurige ooi:

n = het aantal ovulaties,

x = het aantal embryos,

p = de kans dat een willekeurige ovulatie in een embryo resulteert.

Voor ieder van de ovulaties voeren we een indicator variabele in:

$a_m = 1$ als de m -de ovulatie in een embryo resulteert

en 0 als dit niet het geval is, $m=1, \dots, n$.

Nu geldt dus $x = \sum_{m=1}^n a_m$.

De kans p hangt af van de behandeling, leeftijdsklasse en dekperiode, en van het dier zelf.

$$p = p_0 + \epsilon.$$

Hierin is p_0 bepaald door de proeffactoren en is ϵ de stochastische bijdrage van het dier. We nemen aan dat ϵ gemiddelde 0 heeft en een variantie van de vorm:

$$\text{Var}(\epsilon) = \sigma^2 p_0 (1 - p_0).$$

Voor p_0 in de buurt van 0 of 1 is de tussen-dier variatie dus klein en voor p_0 rond 0.5 is de tussen-dier variatie groot.

We veronderstellen dat er sprake kan zijn van een zekere 'competitie' tussen de ovulaties. We voeren daartoe een correlatie ρ in tussen a_m en $a_{m'}$, conditioneel op de dierbijdrage ϵ :

$$\text{cov}(a_m, a_{m'} | \epsilon) = \rho p(1-p) \quad m, m' = 1, \dots, n, \quad m \neq m'.$$

We denken daarbij aan negatieve waarden voor ρ . Aangenomen wordt dat ρ in absolute waarde niet al te groot is. Op z'n minst zal moeten gelden:

$$-1/(n_{\max} - 1) < \rho < 1,$$

waarin n_{\max} het grootste aantal ovulaties is dat kan voorkomen. De correlatie ρ wordt verder constant verondersteld.

De kans p_0 , het 'fixed' gedeelte van p , kan eventueel van n afhangen. Om het model niet te complex te maken is hiervoor een sterk vereenvoudigde voorziening getroffen. Naast de eerder geïntroduceerde factoren behandeling (B, niveau 1 is de proefgroep en niveau 2 de controlegroep), leeftijd (L) en dekperiode (D) voeren we een factor A (van 'aantal') in met 2 niveaus:

$$A = 1 \quad \text{wanneer } n \leq 2 \text{ en } 2 \text{ wanneer } n > 2.$$

De kans p_0 hangt nu van B, L, D en A af.

Voor het modelleren van p_0 gebruiken we de logit-transformatie:

$$\text{logit}(p_0) = \log\left(\frac{p_0}{1-p_0}\right) = \text{algemeen gemiddelde} + \text{som van hoofdeffecten en interacties voor de factoren B, L, D en A.}$$

Om een onevenredig groot aantal parameters te vermijden zijn geen interacties met factor A opgenomen.

Voor de verwachting en variantie van het aantal embryos x geldt nu:

$$\begin{aligned} E(x) &= np_0, \\ \text{Var}(x) &= np_0(1-p_0) \{1 + (n-1)\phi\}, \end{aligned}$$

waarin $\phi = \rho + (1-\rho)\sigma^2$, de (niet-conditionele) correlatie tussen a_m en $a_{m'}$.

De variantie is afgeleid in appendix A2. Merk op dat de correlatie ϕ zowel positief als negatief kan zijn. Deze dispersie parameter is het totaal van overdispersie geïntroduceerd door de dierbijdrage ϵ (tussen-dier variatie) en onderdispersie geïntroduceerd door de (conditionele) 'competitie' correlatie ρ (binnen-dier variatie).

Stel dat ϕ bekend is.

Het verband tussen $E(x)$ en $\text{Var}(x)$ ligt nu vast. Dit maakt het mogelijk om de onbekende parameters in $\text{logit}(p_0)$ te schatten met de Quasi Likelihood methode (QL). QL is beschreven in McCullagh en Nelder (1983). De methode is beknopt weergegeven in appendix A3. De praktische uitvoering verloopt probleemloos met het statistische pakket Genstat (Genstat 5, 1987).

Stel dat de vector van parameters β is gepartitioneerd als $\beta = (\beta_1', \beta_2')'$. Voor een toets van de hypothese $H_0: \beta_2 = 0$, kan de Quasi Likelihood Ratio (QLR) worden gebruikt. Onder H_0 volgt $-2 \log(\text{QLR})$ bij benadering een chi-kwadraat verdeling met als vrijheidsgraden het aantal vrije parameters in β_2 (de reductie in dimensie door het opleggen van de nulhypothese), zie hiervoor McCullagh (1983).

Bij het gebruik van Genstat volgt $-2 \log(\text{QLR})$ uit het verschil van de zogenaamde deviances onder het gerespecteerde model ($\beta_2 = 0$) en het onge-restricteerde model.

Een alternatief is de Quasi Wald-toets gebaseerd op de asymptotische normaliteit van de schatters op de logit-schaal.

Uiteraard is ϕ niet bekend. We gebruiken daarom in het voorafgaande een schatting $\hat{\phi}$. Deze schatting kan worden afgeleid volgens een momenten-methode op basis van de grootheid van Pearson die is gesuggereerd door Williams (1982). De schattingsmethode is iteratief, beschreven in appendix A4 en eenvoudig in Genstat te implementeren. Het kan ook met GLIM, zie hiervoor Williams (1982).

4 DE RESULTATEN

Het grootste model dat is bekeken is:

$$\text{logit}(p_0) = \mu + B + L + D + A + B.L + B.D + L.D + B.L.D.$$

De schatting voor ϕ in dit model is: $\hat{\phi} = -0.0751$. De negatieve uitkomst is een indicatie dat binnen-dier 'competitie' inderdaad een rol zou kunnen spelen. In onderstaande tabel zijn de interacties getoetst m.b.v. de QLR.

Interactie $-2 \log(\text{QLR})$ vrijheidsgraden overschrijdingskans

B.L.D	3.9	3	0.27
B.L.	1.1	3	0.78
B.D.	2.1	1	0.15
L.D	1.3	3	0.73

Voor bijvoorbeeld de twee-factor interactie L.D gebruiken we de deviances van de volgende twee modellen:

$$\mu + B + L + D + A + B.L + B.D + L.D$$

en

$$\mu + B + L + D + A + B.L + B.D$$

Naar aanleiding van de toetsingsresultaten reduceren we het model tot alleen de hoofdeffecten:

$$\text{logit}(p_0) = \mu + B + L + D + A.$$

In dit model bedragen de overschrijdingskansen voor B, L, D en A achtereenvolgens 0.06, 0.14, 1.0 en 0.07. We concluderen dat er geen aantoonbaar verschil bestaat tussen de dekperioden. Kijken we voor de factor leeftijd specifiek naar het contrast tussen de klassen 1 en 4 (de jongste en de oudste dieren) dan is de overschrijdingskans 0.04. Leeftijd is dus van belang.

Zonder factor A in het model draagt de factor behandeling zeer significant bij ($p < 0.01$). Globaal ligt het quotient x/n met $n \geq 1$ voor de behandelde dieren 0.12 lager dan voor de controle dieren. Een rangtoets op de quotiënten x/n (waarin nu ook n als stochastisch wordt beschouwd) geeft een zeer significant resultaat ($p < 0.01$). Deze rangtoets is beknopt weergegeven in appendix A5. Dit is geen bewijs, maar wel een sterke aanwijzing voor de juistheid van de volgende bewering:

Voor de behandelde dieren is de embryonale ontwikkeling significant lager dan voor de controle dieren. Echter, dit is voor een deel terug te voeren tot de veel lagere kans op ontwikkeling die gepaard gaat met hogere aantallen ovulaties.

Wanneer factor A niet in het model wordt opgenomen is er geen aantoonbaar leeftijdseffect. Mogelijk werken hier twee krachten tegen elkaar in:

- (i) een hogere leeftijdsklasse geeft een hoger aantal ovulaties, maar bij een hoger aantal ovulaties is de ontwikkelingskans geringer;
- (ii) een hogere leeftijd geeft bij een gegeven aantal ovulaties een hogere kans op ontwikkeling.

Zonder A in het model heffen (i) en (ii) elkaar grotendeels op. Nemen we A op, dan wordt voor (i) (voor een groot deel) gecorrigeerd en komt (ii) significant te voorschijn.

De parameterschattingen op logit-schaal zijn weergegeven in onderstaande tabel.

Parameter	Schatting	Standaard- afwijking
μ	1.41	0.18
B_1	-0.33	0.18
B_2	0.33	0.18
L_1	-0.55	0.35
L_2	-0.20	0.23
L_3	0.18	0.28
L_4	0.58	0.29
D_1	-0.03	0.15
D_2	0.03	0.15
A_1	0.33	0.18
A_2	-0.33	0.18

Schatten we ϕ opnieuw in het hoofdeffectenmodel dan vinden we $\hat{\phi} = -0.1175$. De schattingen veranderen weinig, de overschrijdingskansen worden iets kleiner en liggen voor A en B net boven de 0.05.

5. SLOTOPMERKINGEN

Er zijn aparte analyses uitgevoerd op het aantal ovulaties en het aantal embryo's per dier, deze komen in dit verhaal verder niet aan de orde. In principe is het mogelijk één model te construeren voor ovulaties en embryo's samen. De omvang van de dataset in combinatie met het discrete karakter van de waarnemingen maakt dit echter weinig aantrekkelijk.

Voor de combinatie van quasi-likelihood en schatting van de dispersie parameter ϕ aan de hand van Pearsons grootheid is in de literatuur nog geen theoretische verantwoording gegeven.

6. REFERENTIES

- Conover, M.J. "Practical non-parametric statistics". 2nd edition, John Wiley & Sons, 1980.
- Genstat 5 Committee, Payne R.W. en anderen. "Genstat 5 Reference Manual", Clarendon Press - Oxford, 1987.
- McCullagh, P. "Quasi-likelihood functions". The Annals of Statistics 11, pp. 59-67, 1983.
- McCullagh, P., Nelder, J.A. "Generalized linear models". Chapman and Hall, 1983.
- Searle, S.R. "Linear Models". John Wiley & Sons, 1971.
- Williams, D.A. "Extra-binomial variation in logistic linear models". Applied Statistics 31, pp. 144-148, 1982.

APPENDICES

- A1 Factoren : Behandeling, Leeftijd, Dekperiode
 Variabelen : n(ovulaties), x(embryos)

B L D n x	B L D n x	B L D n x	B L D n x
1 1 1 2 1	1 2 2 3 2	2 1 1 2 1	2 2 2 1 1
1 1 1 2 2	1 2 2 3 1	2 1 1 2 1	2 2 2 1 1
1 1 1 2 2	1 2 2 3 2	2 1 1 1 1	2 3 1 2 2
1 1 1 1 1	1 2 2 2 1	2 1 1 1 1	2 3 1 3 3
1 1 1 2 1	1 3 1 3 3	2 1 2 1 1	2 3 1 2 1
1 1 1 2 2	1 3 1 3 2	2 1 2 1 1	2 3 1 2 2
1 1 2 2 1	1 3 1 2 1	2 1 2 1 1	2 3 1 3 2
1 1 2 2 1	1 3 1 3 3	2 1 2 1 1	2 3 1 2 2
1 1 2 1 1	1 3 1 2 2	2 1 2 1 1	2 3 2 2 1
1 1 2 2 1	1 3 1 2 1	2 1 2 1 1	2 3 2 2 2
1 1 2 2 1	1 3 2 4 4	2 2 1 2 2	2 3 2 2 1
1 1 2 2 1	1 3 2 4 3	2 2 1 1 1	2 3 2 2 2
1 1 2 1 1	1 3 2 4 1	2 2 1 3 2	2 3 2 2 2
1 2 1 2 1	1 3 2 3 3	2 2 1 2 1	2 3 2 2 2
1 2 1 2 2	1 3 2 3 2	2 2 1 1 1	2 4 1 2 1
1 2 1 2 2	1 3 2 2 2	2 2 1 2 1	2 4 1 2 2
1 2 1 4 2	1 4 1 3 2	2 2 1 2 2	2 4 1 2 2
1 2 1 3 2	1 4 1 2 1	2 2 1 2 1	2 4 1 2 2
1 2 1 2 2	1 4 1 3 2	2 2 1 2 2	2 4 1 2 2
1 2 1 2 2	1 4 1 3 3	2 2 1 3 2	2 4 1 3 3
1 2 1 2 1	1 4 1 2 2	2 2 1 2 2	2 4 1 2 2
1 2 1 2 2	1 4 1 3 3	2 2 2 2 1	2 4 1 2 2
1 2 1 3 1	1 4 2 2 2	2 2 2 2 2	2 4 2 2 2
1 2 2 2 2	1 4 2 4 4	2 2 2 2 2	2 4 2 2 2
1 2 2 2 1	1 4 2 2 2	2 2 2 2 2	2 4 2 2 1
1 2 2 2 2	1 4 2 4 3	2 2 2 2 2	2 4 2 2 2
1 2 2 2 1	1 4 2 2 2	2 2 2 2 2	2 4 2 2 2
1 2 2 3 2	1 4 2 2 2	2 2 2 2 2	2 4 2 3 3
1 2 2 2 2	1 4 2 5 2	2 2 2 1 1	2 4 2 3 3
1 2 2 3 3	2 1 1 1 1	2 2 2 2 2	

$$\begin{aligned}
 {}^{54} A2 \text{ Var}(x) &= \text{Var}(E(x|\epsilon)) + E(\text{Var}(x|\epsilon)) = \\
 &= \text{Var}(np) + E(np(1-p) + n(n-1)p(1-p)) = \\
 &= n^2 \text{Var}(\epsilon) + n(1 + p(n-1))E(p(1-p)) = \\
 &= n^2 \text{Var}(\epsilon) + n(1 + p(n-1))(p_0(1-p_0) - \text{Var}(\epsilon)) = \\
 &= np_0(1-p_0)(1 + \phi(n-1))
 \end{aligned}$$

A3 Stel dat p_* een beginschatting is voor p_0 .

$$\begin{aligned}
 E(x) &= np_0 = np_* + n \left[\frac{\partial p}{\partial \beta} \right]'_{\beta_*} (\beta_0 - \beta_*) \\
 &= np_* + np_*(1-p_*) \underline{b}'(\beta_0 - \beta_*),
 \end{aligned}$$

waarin β_* correspondeert met p_* en β_0 de echte waarde van β is met:

$$\text{logit}(p_0) = \underline{b}'\beta_0.$$

We voeren nu een nieuwe afhankelijke variabele in:

$$Y = (x - np_*)/v_* + \underline{b}'\beta_* \quad \text{met } v_* = np_*(1-p_*).$$

Nu geldt

$$E(Y) = \underline{b}'\beta_0 \quad \text{en} \quad \text{var}(Y) = (wv_*)^{-1},$$

$$\text{waarin } w = (1 + \phi(n-1))^{-1}.$$

Met behulp van een gewogen lineaire regressie voor de variabele Y kunnen we nu een verbeterde schatting voor β_0 afleiden. Deze noemen we weer β_* en we itereren tot het proces convergeert.

A4 Stel dat ϕ_* een beginschatting is voor de echte waarde ϕ_0 van parameter ϕ . Bepaal als in A3 een schatting $\hat{\beta}$ met bijbehorende \hat{p} .

Vat de gewichten w samen in een diagonaalmatrix $W = \text{diag}(w)$, en voer de designmatrix X in waarvan de rijen de vectoren \underline{b} zijn behorende bij de verschillende dieren.

Stel W_* correspondeert met ϕ_* en W_0 met ϕ_0 .

Definieer Pearsons grootheid:

$$X_P^2 = \sum \frac{w_* (x - np)^2}{np(1-p)}.$$

Nu geldt $X_P^2 = (\underline{Y}_* - X\hat{\beta})' W_* V_* (\underline{Y}_* - X\hat{\beta})$,
waarin \underline{Y}_* de variabele Y is geëvalueerd voor $\hat{\beta}$ en ϕ_*
en $V_* = \text{diag}(v_*)$ met $v_* = np(1-p)$.

$$\hat{\beta} = QW_* V_* \underline{Y}_*, \quad \text{waarin } Q = X(X'W_* V_* X)^{-1}X'.$$

Vergeten we nu het stochastische karakter van V_* dan volgt (zie Searle (1971), pag. 55):

$$\begin{aligned} E(X_p)^2 &= \text{spoor}(E(X_p)^2) = E(\text{spoor}(X^2)) = \\ &= \text{spoor}[(I - QW_*V_*)'W_*V_*(I - QW_*V_*)(W_0V_*)^{-1}] + \\ &\quad \beta_0'X'(I - QW_*V_*)'W_*V_*(I - QW_*V_*)X\beta_0. \end{aligned}$$

Nu geldt $W_*V_*(I - QW_*V_*)X = 0$, dus

$$\begin{aligned} E(X_p^2) &= \text{spoor}[(I - QW_*V_*)'W_*V_*(I - QW_*V_*)(W_*V_*)^{-1}W_*W_0^{-1}] \\ &= \text{spoor}[(I - V_*W_*Q)W_*W_0^{-1}] = \\ &= \Sigma w_*(1 - v_*w_*q)/w_0 \end{aligned}$$

met q de diagonaal van Q .

Los nu een nieuwe waarde voor ϕ op uit:

$$X_p^2 = \Sigma w_*(1 - v_*w_*q)(1 + (n-1)\phi).$$

Itereer tot convergentie is bereikt. In dat geval geldt:

$$\begin{aligned} X_p^2 &= \text{spoor}[I - V_*W_*Q] = \text{rang}[I - V_*W_*Q] = \\ &= N - \text{rang}(X), \text{ daer } V_*W_*Q \text{ idempotent is, } N \text{ is het aantal dieren.} \end{aligned}$$

Het proces kan worden gestart met $\phi_* = 0$.

A5 Bepaal per dier x/n met $n \geq 1$.

Geef binnen iedere combinatie van factoren L en D rangnummers aan de quotienten. Laat T voor een combinatie de som van de rangnummers van de behandelde dieren zijn (dit is de toetsingsgrootte van de toets van Wilcoxon). Onder de hypothese H_0 : Fecundin heeft geen effect, geldt:

$$E(T) = a(b+1)/2, \quad \text{Var}(T) = \frac{a(b-a)}{b(b-1)} \Sigma R^2 - \frac{a(b-a)(b+1)^2}{4(b-1)},$$

waarin a het aantal behandelde dieren is, b het totaal aantal dieren en ΣR^2 de som van de kwadraten van de rangnummers van de behandelde dieren voor de betreffende combinatie, zie Conover (1980). Refereer $(\Sigma T - \Sigma E(T))/(\Sigma \text{Var}(T))$ aan de standaard normale verdeling. Aangenomen is dat een eventueel effect van Fecundin consequent positief of negatief is voor de combinaties van L en D . Dit is redelijk te controleren met behulp van de toets van Wilcoxon uitgevoerd per combinatie van L en D . Alleen voor dekperiode 1 en leeftijdsklasse 1 bleek $T < E(T)$ (niet significant), voor de overige combinaties was $T > E(T)$.