

1. Inleiding

Dit is een verslag van lopend onderzoek naar de methodologie van het analyseren van kleinsteekproeven voornamelijk van discrete data. We bekijken en vergelijken een aantal technieken voor het bepalen van steekproevenverdelingen, en we maken software hiervoor, waarbij gebruik wordt gemaakt van efficiëntie doelen. We zoeken naar regels die de toepassing van deze technieken en de beperkingen van elk van deze technieken aangeven. Tenslotte noemen we een aantal methodologische problemen die specifiek zijn voor kleine steekproevenverdelingen, die bij toenemende steekproefomvang vanzelf verdwijnen en die mede daardoor in het algemeen weinig aandacht krijgen.

Hoewel deze problemen de zelfde oplossingsmethoden in zeer verschillende situaties tegenkomen, zijn we er (nog?) niet in geslaagd een algemene probleemstelling te formuleren waarvan de speciale gevallen voortvloeien. We beperken ons daarom tot een opsomming van de drie belangrijkste situaties waarop ons onderzoek zich richt.

1.1 — 2x2- en grotere rxc-tabellen:

Chances op onafhankelijkheid of homogeniteit. Permutatie- of exacte tests tegen geordende alternatieven, zoals de correlatieratio, Kruskal-Wallis, Kendall's tau, Spearman's R_s en de gewone correlatiecoëfficiënt. Indien één variabele (correlatieratio, Kruskal-Wallis) of beide variabelen (tau, H_s , of r) van ordinaal of interval-niveau zijn, dan zijn er dus veel knopen. Anders zouden we de tests niet in een rxc-tabel presenteren.

Kleine-steekproevenverdelingen en betrouwbaarheidsintervallen voor associatiematen.

1.2 — Methoden gebaseerd op de permutatieverdeling van rangnummers of van de oorspronkelijke scores (met weinig of geen knopen).

Hierbij denken we vooral aan non-parametrische toetsen.

1.3 — Meerdere dimensionale kruistabellen: loglineaire en logit-analyse.

χ^2 -goodness-of-fit toetsen.

Verdelingen van de schakelaars.

2. Het bepalen van steekproevenverdelingen

Voor het eventueel bij benadering bepalen van steekproevenverdelingen van een stochast, zeg T , kunnen de volgende vijf methoden nuttig zijn:

2.1 — Tabellen.

Deze zijn alleen beschikbaar voor eenvoudige situaties met weinig parameters en (erg) kleine aantallen waarnemingen. Zouden bijvoorbeeld knopen optreden, of gecensureerde data, of stratificatie van de steekproef in de tabellen ontbreken, in het algemeen in de steek.

2.2 — Het gebruik van asymptotische methoden.

De meeste goodness-of-fit toetsen zijn gebaseerd op de likelihood ratio toets. Af asymptotisch gelijke grootheden zoals $\chi^2 = \sum (\text{obs} - \text{exp})^2 / \text{exp}$. Voor loglineaire en logit-modellen is de maximum-likelihood methode de universeel gebruikte schattingsmethode. Ook bij permutatietoetsen bij geknoopte waarnemingen wordt meestal gebruik gemaakt van de asymptotische verdeling.

Het grote probleem met de meeste asymptotische methoden is dat weinig bekend is over hun nauwkeurigheid bij kleine steekproeven: wolkomen of bijna "voldoende groot"? Mede doordat in veel toetsen de steekproefomvang van één of twee cijfers ruimschoot kan variëren, is de gebruikte asymptotische methoden vaak niet bij voldoende kleine steekproeven voldoende nauwkeurig. Dit gaat, dan vaak goed, maar nogmaals: er zijn weinig regels die ons vertellen of het goed gaat.

2.3 — Het aftellen van alle punten uit de uitkomstenruimte.

De voornaamste beperking van deze methode is dat het al gauw veel werk is. Zo is het aantal permutaties op niveau de orde $(n/e)^n$, en voor $n=10$ is dat ruim 30 miljoen. Het aantal rxc-tabellen met gegeven marginaal is $n!$, en voor $n=10$ is dat ruim 30 miljoen (maar niet exponentieel, zoals vaak gedacht wordt). Het is dus geen enkel bezwaar bij

Verbeek, A. & Kroonenberg, P.M. (1982)

Kleine steekproevenverdelingen (Samenvatting 'lezing
Statistische dag VVS).

VVS Bulletin, 15 (5), 18-19, 22-23.

KLEINE-STEELPROEVENVERDELINGEN

Voordracht voor de Statistische Dag, U april 1982, Erasmus Universiteit Rotterdam

Albert Verbeek, Biologisch Instituut, Rijksuniversiteit Utrecht

Pieter H. Kroonenberg, vakgroep WEP, subfaculteit PAW, Rijksuniversiteit Leiden

Inhoud:

- 1 Inleiding.
Belangrijke onderzoeksituaties:
 - 1.1 2×2 - of $r \times c$ -tabel "I" "n"
 - 1.2 permutatieverdelingen
 - 1.3 Loglineaire en logit-analyse
- 2 Het bepalen van steekproevenverdelingen
 - 2.1 Tabellen
 - 2.2 Asymptotische methoden
 - 2.3 Het aftellen van de uitkomstenruimte ("exacte" methoden)
 - 2.4 Monte-Carlo en bootstrapmethoden
 - 2.5 Karakteristieke functies
- 3 Toetsproblemen
 - 3.1 Toetsing van de toetsingsgrootte; onderscheidingsvermogen
 - 3.2 Wat is een "kleine steekproef"?
 - 3.3 Discreetheit en symmetrie in spelbedervers
 - 3.4 $IK(\text{STATISTIC}, \text{OBSERVEDVALUE}) \dots$ met STATISTIC en OBSERVEDVALUE reëel
 - 3.5 Conditioneren en randomiseren
 - 3.6 Toevallige nulcellen

1. Inleiding

IMT is een verslag van lopend onderzoek naar de methodologie van het analyseren van kleine steekproeven voornamelijk vnn discrete data. We bekijken en vergelijken ook aantal technieken voor het bepalen vnn steekproevenverdelingen, en we maken software hiervoor, waarbij gebruiksvriendelijkheid, ruime beschikbaarheid en efficiëntie doelen zijn. We zoeken naar n... de toepasbaarheid en de beperkingen van elk vnn do7o technieken aangeven. Tenslotte noemen we een aantal methodologische problemen die specifiek zijn voor kleine steekproevenverdelingen, die bij toenemende steekproefomvang vanzelf verdwijnen en die mede door dit hot, n]gemeen weinig aandacht krijgen.

Hoewel «o dezelfde problemen op dezelfde oplossingen in zeer verschillende situaties tegenkomen, zijn we or (nog??) niet in geslaagd één algemene probleemstelling te formuleren waarvan die situaties speciale gevallen zijn. We beperken ons daarom tot een opsomming vnn do drie belangrijkste situaties waarop ons onderzoek zich richt.

1.1 — 2×2 - Pn r... x-c-tabellen:

χ^2 -toetsen op onafhankelijkheid of homogeniteit. Permutatietoetsen tegen geordende alternatieven, zoals do correlatiecoëfficiënt, Kruskal-Wallis, Kendall's tau, Spearman's R_s en do gewone correlatiecoëfficiënt r . Indien één variabele (correlatiecoëfficiënt, Kruskal-Wallis) of beide variabelen (tau, R_s or r) van ordinaal of interval-niveau zijn, dan zijn or dus veel knopen. Anders zouden we de data niet in een x-c-tabel presenteren.

Kleine steekproevenverdelingen en betrouwbaarheidsintervallen voor associatie-maten.

1.2 — Methoden gebaseerd op do permutatieverdeling van rangnummers of vnn do oorspronkelijke scores (met weinig of geen knopen).

Hierbij denken we vooral aan non-parametrische toetsen.

1.3 — Meerdimensionale kruistabellen: loglineaire or logit-analyse.

χ^2 -goodness-of-fit toetsen.
Verdelingen vnn d" schalliers.

2 Het bepalen van steekproevenverdelingen

Voor het eventueel bij benadering bepalen vnn steekproevenverdelingen van een stochast, zeg T , kunnen do volgende vijf methoden nuttig zijn:

2.1 — Tabellen.

Deze zijn alleen beschikbaar voor eenvoudige situaties met weinig parameters or (erg) kleine aantallen waarnemingen. Zodra bijvoorbeeld knopen optreden, of gecensureerde data, or stratificatie vnn de steekproef laten tabellen or in het algemeen in do steek.

2.2 — Het gebruik vnn asymptotische methoden.

De meeste goodness-of-fit toetsen zijn gebaseerd op de likelihood-ratiotoets or asymptotisch gelijke grootheden zoals Pearson's $\chi^2 = \sum (obs - exp)^2 / exp$. Voor loglineaire or logit-modellen is do maximum-likelihoodmethode de meest gebruikte schattingsmethode. Ook bij permutatietoetsen bij geknoopte waarnemingen wordt meestal gebruik gemaakt vnn d" asymptotische verdeling.

Het grote probleem met d" meeste asymptotische methoden is dat weinig bekend is over hun nauwkeurigheid bij eindige steekproeven: welke steekproeven zijn "voldoende groot"? Mede doordat il. veel toepassingen or nauwkeurigheid van één of twee cijfers ruimschoots voldoende is, zijn do gebruikelijke asymptotische methoden vaak al bij verrassend kleine steekproeven voldoende nauwkeurig. Dit gaat dus vaak goed, maar nogmaals: or zijn weinig regels of richtlijnen or het goed gaat.

2.3 — Het aftellen vnn alle punten uit do uitkomstenruimte.

De voornaamste beperking vnn deze methode is dnt. bot, ni gauw voor werk is. Zo is het aantal permutaties nivan de orde $(n/e)^n$, or voor $n = 10$ is ni ni ruim 30 miljoen. Met, aantal x-c-tabellen met gegeven marginals neemt toe als n^2 voor $n \rightarrow \infty$ (maar niet exponentieel, zoals vnnk gedacht wordt). Dit is dus, voor kleine n , bezwaar bij

2x2-tabellen (waarbij v_g I), ondanks dat SPSS vanaf $n=21$ zeer dubieuze asymptotische benaderingen geeft. Mot, enige behendigheid en mot de snelheid van hodendaags rekentuig is het aftellen van de uitkomstenruimte in een aantal interessante gevallen toepasbaar. Een andere aardige toepassing is het gebruik van deze methode hij het statistiekondorwijs.

2.4 — Monte-Carlo en bootstraphethoden.

Als we mbv. pseudorandom getallen enkele honderden tot enkele duizenden trekkingen genereren nlt. een bepaalde verdeling, en wo berekenen voor elke trekking de waarde vnn T (do stochast waar we in geïnteresseerd zijn), dan krijgen wo een tamelijk nauwkeurige schatting vnn do verdeling van T. Als do verdeling wnn wo telkens uit trekken bekend in (zoals hij een enkelvoudige H_0) spreken we van Monte-Carlo; als we deze verdeling eerst mbv. do data moeten schatten, noemen we het pon bootstraphethoden. Inkoop is dit in bol, algemeen niet: er moot. In de orde vnn 1000 keer dezelfde schatter berekend worden. Is veel gevallen is dit duur maar niet onmogelijk.

2.5 — Karakteristieke functies.

In sommige gevallen is mogelijk een recursieve betrekking aan te geven waarin do karakteristieke functie vnn T (dzw $E \exp(iTt)$) voor n waarnemingen nltgedrukt wordt in dio voor $(n-1)$ waarnemingen. Deze methode in voorgesteld door Paganofit Tritchler, o.a. voor Wilcoxon's één- en tweesteekproeventoetsen, waarbij de omvang van do berekeningen bij toenemende n slechts evenredig zijn mot, n . Voor het berekenen vnn kannen uit do karakteristieke functie kan men do Fast Fourier Transform gebruiken, waarbij do rekentijd niet afhangt vnn n (maar wel. van de gewenste precisie).

3 Specifieke kleine-steekproevenproblemen

3.1 Keus van d toetsingsgroottheid bijv. bij toetsen op onafhankelijkheid in een rxc-tabel.

Do meest gebruikt groottheden zijn: Pearson's χ^2 , do $-2 \log$ likelihood ratio, Freeman-Tukey en $-2 \log p(\text{tabel} | H_0)$. Deze groottheden zijn asymptotisch gelijk, "n dur: geeft de asymptotische theorie weinig houvast voor non keuze bij kleine steekproeven. Vnak wordt dnn een tonts gekozen waarvan de verdeling het best benaderd wordt door de χ^2 -verdeling (waarschijnlijk kiest mon dan χ^2). Maar dat is natuurlijk eigenlijk oen tweederangs argument. (Hoewel, wat heb je aan een betere toets waarvan je do verdeling niet kunt. bepalen?)

Ton belangrijke vraag is ook: hoe groot moot, do steekproef minstens zijn om toetsen überhaupt zinvol te maken? In oen kruistabel met, alleen nullen en enen is het onderscheidingsvermogen vnn bovengenoemde toetsen zo gering (denken we), dat jo not, 7,0 goed at random met kans α zou kunnen verwerpen.

Het lijkt een goed idee om in standaardprogrammatuur 7,0 te maken dat ook enige zinvolle informatie gegeven worden over het onderscheidingsvermogen van do gekozen toets(en). Omdat de alternatieve hypothese meerdimensionaal is, is het echter niet zo makkelijk om automatisch enkele relevante verdelingen uit de alternatieve hypothese te kiezen.

Veelvuldig komt mon do opvatting tegen, dnt voor het berekenen van exacte overschrijdingskansen de punten u i. de uitkomstenruimte geordend zouden moeten worden nnnr afnemende kansen. In feito neemt mon dan $-\log$ kans als toetsingsgroottheid. De voornlanders vinden dit 7,0 vanzelfsprekend, dnt, 7,0 do nltornatieve ven niet eens noemen. O. i. berust dit op "on misvatting, on we 7,0 ion geen reden om $-\log$ kans moor of minder aantrekkelijk te vinden dan de andere toetsingsgroottheden. Dat het makkelijker te berekenen in, is woor een tweederangsargument.

3.2 Cochran heeft in jaar geleden de volgende vuistregel gegeven voor do χ^2 -benadering voor Pearson's χ^2 in rxc-tabellen onder onafhankelijkheid, onder do voorwaarde $(r-1)(c-1) > 1$.

Indien allo verwachte waarden groter zijn dnn 1 on minstens 80% is groter dnn 5, en men neemt bot χ^2 -95%-punt als drempelwaarde, dan ligt het werkelijke niveau vnn deze toets tussen de 3 en 7%.

In plaats hiervan stellen wij do volgende vuistregel voor, din niet scherper is, mnr wol ruimer toepasbaar:

Indien de rijmarginale frequenties onderling verschillend zijn, en de kolommarginale frequenties $\delta\delta k$, en men neemt het χ^2 -95%-punt als drempelwaarde, dan ligt het werkelijk niveau van α iets minder ver van de (ideale) 0,05 af dan

$$\frac{0,05}{r \cdot c} \cdot \frac{1}{\exp}$$

- 3.3 In veel gevallen blijkt de asymptotische verdeling vooral van de echte verdeling te verschillen doordat de laatste discreet is. In situaties waarin deze discreetheid extra geprononceerd is, is bijvoorbeeld een rxc-tabel met gelijke marginales (vergelijk de conditie in de eerste regel van deze pagina). De gangbare opvatting dat de asymptotische benaderingen vooral in symmetrische situaties goed toepasbaar zijn, is dus niet algemeen waar.
- 3.4 Indien de toetsingsgrootte geen geheel getal is, en er zijn verschillende punten in de uitkomstenruimte met dezelfde of vrijwel dezelfde waarde van de toetsingsgrootte, dan kunnen afrondfouten in de berekening de volgorde verstoren. Het gevolg kan zijn dat verschillende computers met hetzelfde programma voor dezelfde data een wezenlijk verschillende overschrijdingskans berekenen.

Soms kunnen we dit probleem omzeilen door de toetsingsgrootte naar een geheelwaardige variabele i_r te transformeren (zonder in overflow-problemen te geraken). Soms ook kunnen we een α aangeven zó dat waarden die minder dan α verschillen gelijk zijn. En als dat allemaal niet lukt, dan kunnen we altijd een ondergrens en een bovengrens voor het gevraagde antwoord berekenen. (De meeste programma's geven alleen een bovengrens voor de p-waarde, overigens zonder α te vermelden dat, dit, alleen op bovengrens is...)

3.5 Conditioneren en randomiseren.

De vraag of je bij het toetsen op onafhankelijkheid in een rxc-tabel moet conditioneren op de marginales, en de vraag of je mag of moet randomiseren leveren lovende discussies op (bijv. Sniijders, 1980).

Hot is typisch een probleem dat bij toenemende steekproefomvang verdwijnt. Omdat wel of niet conditioneren asymptotisch geen verschil maakt, terwijl het hot effect van randomiseren bij afnemende discreetheid verdwijnt.

- 3.6 Bij het analyseren van log-linear modellen kan het hot gebeuren dat de enkoortevalle-nulcellen het schatten van bepaalde parameters onmogelijk maken. Het kloint steekproeven kan het zelfs gebeuren dat je helemaal geen modeltoets kunt uitvoeren. Hoewel de meeste handboeken deze situatie wel behandelen, vermelden ze nooit of de door hun beschreven toets niveau n heeft, conditioneel op een uitkomst waarbij getoetst kan worden, of onconditioneel. Hoewel het hier om een tamelijk extreme situatie gaat, is het toch opvallend hoe weinig gebruikers van deze methoden hierover nagedacht hebben. Het is woeer typisch zo'n probleem waarvan de kans op optreden al snel heel klein wordt bij toenemende steekproefomvang.

Referenties

De programmabeschrijving van FISHER, in KM 2(1981) p. 66-86 bevat een uitvoerige lijst referenties, maar niet de volgende:

H. Pagano & K. Taylor Halvorsen - An algorithm for finding the exact significance levels of rxc contingency tables. JASA 76 (1981) 931-934

M. Pagano & D. Tritchler - On obtaining permutation distributions in polynomial time. Prepublicatie Harvard/Sidney Farber Cancer Institute. (1980)

W.M. Mitchell - An efficient method of generating random rxc tables with given row and column totals. Applied Statistics 30 (1981) AS159

De discussie over wel of niet conditioneren met Tom Sniijders: zie VWS-Bulletins februari en maart (1980).