

Een blik op Databases, Types and the Relational Model

Twijfels over logische correctheid

Maurice Gittens

In dit artikel wordt beargumenteerd dat de nieuwste editie van The Third Manifesto 'Databases, Types and The Relational Model' van Date en Darwen vanuit sommige perspectieven een regressie is ten opzichte van het relationele model, zoals deze door Codd is geïntroduceerd en uitgewerkt.

Ook in de nieuwste editie van The Third Manifesto (TTM) wordt de grondregel 'Alle logische verschillen zijn grote verschillen' en zijn uitvloeisel 'Alle logische fouten zijn grote fouten' als maatstaf voor de dissertatie gepresenteerd. Om deze reden heb ik besloten om na te gaan of deze maatstaf door Date en Darwen in deze nieuwste editie wordt gehaald. Evident is al gauw dat, significant meer dan in de vorige editie, deze maatstaf wordt weerspiegeld in de kwaliteit van de definities, proscripties en prescripties in dit boek. Opvallend is ook dat er in deze derde editie van TTM met geen woord gerept wordt over de zogenoemde grote blunders die in de tweede editie van het boek zo prominent aanwezig zijn. Dit artikel presenteert een aantal significante, overblijvende punten waardoor TTM het predicaat 'logisch consistent' niet verdient en misschien zelfs als een regressie ten opzichte van het werk van Codd gezien kan worden.

De TTM vervangt operatoren uit de algebra van Codd met relaties zonder hierbij al te secuur te werk te gaan. Het evidente punt dat commutatieve operatoren als relaties te zien zijn, wordt aannemelijk gemaakt en aangedragen als reden voor de aanpassing van de algebra van Codd. Date and Darwen hebben hierbij nagelaten om:

- aannemelijk te maken dat de operatoren, die vervangen zijn, ook commutatief zijn, of;
- aannemelijk te maken dat niet-commutatieve operatoren ook als relaties te zien zijn.

Date en Darwen hebben niet aangetoond dat het mogelijk is om binnen de door hun opgelegde prescripties en proscripties de bovenstaande punten aannemelijk te maken. Daardoor kan gesteld worden dat Date en Darwen de maatstaf van logische correctheid in deze niet hebben gehaald.

Semantische Compositie

De logicus Frege introduceerde in de logica het idee van de semantische compositie. Vereenvoudigd houdt dit in dat het voor alle expressies in een logische taal mogelijk dient te zijn om een

variabele te vervangen door de waarde van de variabele, zonder dat daarbij de betekenis van de betreffende expressies verandert. Date en Darwen stellen de logica als basis van hun betoog in TTM. Desondanks voldoet de TTM niet aan dit belangrijke principe. Dit blijkt als we bemerken dat volgens TTM:

- het type van een relatie door de header van de relatie wordt bepaald;
- de header van relaties niets zegt over de candidate keys die voor de betreffende relaties gedefinieerd zijn.

Als gevolg hiervan is het mogelijk dat twee relvars r_1 en r_2 , hoewel van hetzelfde type, mogelijk niet aan elkaar toe te wijzen zijn omdat de voor r_1 en r_2 gedefinieerde candidate keys mogelijk verschillend zijn. Zoals hieronder getoond wordt, kan de waarde van r_2 niet aan r_1 worden toegekend terwijl deze relaties van hetzelfde type zijn, dit omdat er een 'unique constraint violation' op de attribute CODE bij de toekenning optreedt (candidate keys worden in de tabellen met hoofdletters aangemerkt).

Relvar r_1		Relvar r_2	
NAME	CODE	NAME	code
Hugh	20	Hugh	20
Chris	25	Chris	20

Als gevolg van het bovenstaande is het formalisme dat in de TTM is gepresenteerd niet referentieel transparant. Dit vertegenwoordigt een logische fout. Opmerkelijk is ook dat TTM voorschift 21 en het door Date en Darwen hoog gedragen principe van conceptuele integriteit overtreden worden door deze fout. Het blijkt namelijk dat wanneer een relvar V en een relatiewaarde v , van hetzelfde type T , aan elkaar worden toegekend volgens de toekenning $V:=v$, dat niet in alle gevallen geldt dat de expressie $V=v$ waar zal zijn. Deze fout is ook terug te vinden als we in ogen-schouw nemen dat de TTM candidate key constraints en foreign key constraints toelaat op relvars, maar niet op relatiewaarden. Zoals ook later in dit artikel blijkt ondermijnt dit feit de waarde van de door de TTM geïntroduceerde 'relation valued attributes' op serieuze wijze.

Geen semantische integriteit bij toekenningen

Date en Darwen stellen dat update-, delete- en insert-operaties alle speciale gevallen van de toekenningoperator zijn. In deze sectie wordt geïllustreerd dat de toekenningoperator zoals deze door TTM gedefinieerd is, een regressie is ten opzichte van Codd (RM/T), omdat het de semantische integriteit van databases ondermijnt. Beschouw eens de volgende relatie met als naam Auteur:

SURNAME	first name
Date	Chris
Darwen	Hugh

Conform de voorschriften van TTM leest de relvar predicaat (zie pagina 29 TTM) voor deze relatie:

Auteur met achternaam 'Date' heeft als voornaam 'Chris'.

Auteur met achternaam 'Darwen' heeft als voornaam 'Hugh'.

Beschouw het volgende toekenningstatement dat de sleutelwaarden van bestemde tuples verwisseld.

```
update Author where surname = "Date"
    { surname = "Darwen" },
update Author where surname = "Darwen"
    { surname = "Date" };
```

Zoals onder meer op pagina 179 van TTM kan worden vastgesteld, is dit statement volgens TTM voorschriften geldig, mede omdat integriteitcontroles pas na het volledige statement worden uitgevoerd. De voorbeeldrelatie wordt dus:

SURNAME	first name
Darwen	Chris
Date	Hugh

De corresponderende relvar predicaat leest dan:

Auteur met achternaam 'Darwen' heeft als voornaam 'Chris'.

Auteur met achternaam 'Date' heeft als voornaam 'Hugh'.

Gegeven de kennis dat er slechts een statement is doorgevoerd op de voorbeeldrelatie en op de waarden van deze relatie voor en na dit statement: wat zou dan de conclusie zijn van een forensische applicatie die belast is met de vraag 'Wat is er veranderd door

dit statement'? Het voor de hand liggende en evident onjuiste antwoord op deze vraag is:

- de voornaam van de auteur met achternaam 'Date' is vervangen door 'Hugh', en;
- de voornaam van de auteur met de achternaam 'Darwen' is vervangen door 'Chris'.

Dit antwoord is niet alleen onjuist. Het is mogelijk nog erger dat we niet op basis van de beschikbare informatie kunnen achterhalen welke wijziging er juist is doorgevoerd.

Hierdoor is na één samengesteld toekenningstatement de semantische integriteit van TTM databases in twijfel te trekken. Heeft u hier moeite mee, dan verzoek ik u het bovenstaande voorbeeld te doordenken in de context van duizenden financiële transacties op bankrekeningen per dag. In deze context zijn individuele tuples aan te wijzen die met bankrekening- en saldogegevens van rechtspersonen corresponderen. Door het uitvoeren van deze denkexerctie zult u begrijpen dat het noodzakelijk is om wijzigingen op individuele tuples in toekenningstatements te kunnen traceren. Dit kan natuurlijk op verschillende manieren, bijvoorbeeld:

- door terug te grijpen op surrogaten (niet muteerbare OID's) zoals deze door Codd zijn geïntroduceerd;
- door bij assignment statements aan te geven welke candidate key waarden niet mogen wijzigen.

Wie als informatiebeveiliging of business rules professional deze kwestie beschouwt, zal gauw erkennen dat tuple-identiteit een vereiste is om wijzigingen op individuele tuples te kunnen toetsen op hun conformering aan informatiebeveiligings-constraints en business rules in het algemeen. Het stokpaardje van Date en Darwen dat tuples in relaties geen eigen identiteit mogen hebben, is om deze en ook andere redenen (zie mijn eerder in DB/M verschenen artikelen over hogere orde concepten) wat mij betreft aan herziening toe.

Ondermijnende issues met relation valued attributes

TTM introduceert een eigen variant van het idee van relation valued attributes. De wijze waarop Date en Darwen dit gedaan hebben verdient echter, vanuit het perspectief van de logische consistentie, niet de schoonheidsprijs. Dit baseer ik op de volgende gronden:

- de algebra van TTM definieert geen operatoren die relation valued attributes ter manipulatie isoleren. Dit heeft tot gevolg dat vanuit het perspectief van de TTM algebra 'relation valued attributes' niet te onderscheiden zijn van types zoals strings, integers of XML documenten. Wat is dan vanuit dit perspectief de meerwaarde van deze 'relaties'?;
- tevens geldt dat TTM niet voorziet in faciliteiten waarmee relation valued attributes van eigen candidate keys voorzien kunnen worden; dit impliceert dat: a. er minder uniciteit constraints op deze 'relaties' gedefinieerd kunnen worden; b. dat

deze 'relaties' niet de rol van parent in 'foreign key constraints' kunnen aannemen, omdat hiervoor candidate keys vereist zijn. Om het nog schrijnender te maken geeft de TTM niet één voordeel van de inzet van deze relation valued attributes, terwijl de nadelen evident zijn. Date en Darwen stellen zelfs niet te weten wanneer deze 'relaties' het beste van pas komen. Op pagina 379 kunt u dit in de context nalezen. Daarom de vraag: Waaruit blijkt dat deze relation valued attributes niet een doel op zich zijn? Vertegenwoordigen deze 'relaties' in het licht van het voorgaande in logische zin een verbetering ten opzichte van Codd? Ik stel van niet.

Een nauwe appreciatie van het onbekende

Codd, samen met veel praktikanten, onderkent de dagelijkse realiteit van het onbekende. In zijn model voor databases heeft hij dit concept geaccommodeerd in de vorm van de ondersteuning van nulls. Date en Darwen verwerpen al jaren dit concept in hun incarnatie van het relationele model. De argumentatie voor deze verwerping van nulls is veelal SQL-specifiek. Ook deze nieuwste versie van TTM biedt geen wezenlijke onderbouwing voor dit standpunt van Date en Darwen. Dit is jammer omdat het goed was geweest om hun redenen voor de verwerping van het concept van een null te zien in een context die niet SQL gerelateerd is. Het feit dat SQL niet handig met nulls omgaat, is namelijk niet een logisch geldige reden voor de verwerping van nulls. Ook de

Er is absoluut niets verkeers met het intelligent ondersteunen van het concept van het onbekende

aanname dat ondersteuning van nulls de ondersteuning van driewaardige logica impliceert is niet zonder meer juist. Waarom zou men niet een taal kunnen ontwerpen of een algebra kunnen definiëren die het concept van het onbekende ondersteunt, zonder dat daarbij driewaardige logica noodzakelijk is? Beschouw in dit licht bijvoorbeeld de volgende propositie die correspondeert met een relvar predicaat uit TTM, zie pagina 29.

```
Supplier S1 is under contract, is named Smith,  
has status 20, and is located in city London.
```

Als we aannemen dat het attribuut 'city' optioneel is en dat de waarde van dit attribuut onbekend is, kunnen we uit ten minste drie formuleringen kiezen voor het relvar predicaat:

1. Supplier S1 is under contract, is named Smith, has status 20, and is located in an unknown city;
2. De relvar verwerpen als niet relationeel;
3. Supplier S1 is under contract, is named Smith, has status 20.

Talen zoals SQL kiezen voor de eerste optie. Date en Darwen kiezen voor de tweede optie. Kan er niet een taal worden ontworpen die uitgaat van de derde optie? Beschouw eens de volgende relatie R met een attribuut 'city' dat onbekend mag zijn. Aan hen die stellen dat dit geen relatie is, voeg ik toe dat dit volgens Codd wel een relatie is.

SUPPLIERID	name	city (mogelijk null)
S1	John	
S2	Jane	Paris
S3	Judy	

Als een taal waarborgt dat operatoren altijd relaties zonder onbekende attributen retourneren, dan is het niet erg dat in de onderliggende database attributen die onbekend zijn bestaan.

Dit wil bijvoorbeeld zeggen dat een statement als:

```
select * from R
```

een foutmelding zou produceren, terwijl de statements:

```
select SUPPLIERID, name from R
```

en

```
select * from R where city="Paris"
```

zonder fouten zouden worden geëxecuteerd, aangezien de resulterende relaties niet naar het onbekende verwijzen. Essentieel is de onderkenning dat ondersteuning voor het concept van het onbekende in databases niet een verplichting is om null-waarden te introduceren in query-talen. Het miskennen van het logische verschil tussen de accommodatie van het onbekende en de ondersteuning van null-waarden en driewaardige logica is dan ook een logische fout. Naar mijn mening is er absoluut niets verkeers met het intelligent ondersteunen van het concept van het onbekende.

Niet in het minst omdat het mogelijk is om dit te doen zonder dat hierbij driewaardige logica een noodzaak wordt. Heeft TTM aannemelijk gemaakt dat het onderwerp nulls goed is doordacht en dat dit onderwerp zonder vooringenomenheid is overwogen? Kan gegeven deze feiten worden gesteld dat Date en Darwen pertinente en valide argumentatie hebben aangedragen voor de verwerping van het concept van null? Ik stel van niet.

Conclusies

Volgens Date en Darwen geldt: 'Alle logische verschillen zijn grote verschillen' en zijn uitvloeisel 'Alle logische fouten zijn grote fouten'. Dit artikel heeft een aantal logische fouten in TTM aangewezen, telkens op punten waarop Date en Darwen een andere richting dan Codd hebben gekozen. Deze punten zijn:

- twijfel over de logische correctheid van de invalshoek 'Operatoren zijn relaties', zoals deze door Date en Darwen is

gepresenteerd, is verantwoord omdat Date en Darwen in hun behandeling van dit thema *niet-commutatieve* operatoren achterwege hebben gelaten;

- ondanks dat Date en Darwen zich logische correctheid toe-eigenen, geldt dat het formele systeem dat TTM introduceert voorbij gaat aan Frege's principe van 'semantische compositie'. Dit is een logische fout, niet in het minst omdat als gevolg van deze fout Date en Darwen in overtreding zijn van hun eigen RM voorschrift 21 en ook het door hen gebezigd principe van conceptuele integriteit.
- Relation Valued Attributes zoals die door TTM geïntroduceerd zijn, zijn vanuit het perspectief van de TTM algebra niet te onderscheiden van andere datatypen zoals strings en XML documenten. Tevens bieden deze attributen minder mogelijkheden voor het faciliteren van uniciteit constraints dan alternatieven. Het bestaansrecht van deze attributen is dan ook in twijfel te trekken, zeker als we in ogenschouw nemen dat als gevolg van dit feit TTM RM voorschrift 26 wordt ondermijnd;
- de semantische integriteit van TTM databases is op basis van de TTM definitie van toekenning niet verzekerd. Dit is gerelateerd aan het feit dat Date en Darwen, anders dan Codd, het concept van tuple identity hebben verworpen;
- anders dan Codd verwerpen Date en Darwen het concept van null. TTM biedt geen bewijs dat Date en Darwen met een open blik naar de ondersteuning van nulls hebben gekeken. Ook is

op geen enkele wijze aannemelijk gemaakt dat nulls op logische gronden niet wenselijk zouden zijn. Hierdoor is de verwerping van nulls door Date en Darwen als dogmatisch te zien.

Het is opmerkelijk te noemen dat de bovenstaande punten juist die punten zijn waarop het relationele model van Date en Darwen in wezenlijke zin verschilt van dat van Codd. Om deze reden is, mijn inziens, de vraag of Date en Darwen ten opzichte van Codd een regressief pad hebben ingeslagen terecht.

Maurice Gittens is zelfstandig IT-consultant.

Online archief Database Magazine

Database Magazine-lezer opgelet! Artikelen over onderwerpen als Datawarehousing, SQL, ETL, Business Intelligence, Relationale databases, modellering en nog veel meer vindt u in het Online Archief van Array Publications. Vaktijdschriften als Storage Magazine, Database Magazine, IT Service Magazine, Java Magazine en ons Oracle vakblad Optimize hebben hun artikelenarchief online gezet. Met een Google-achtige zoekstructuur vindt u snel wat u zoekt op www.dbm.nl



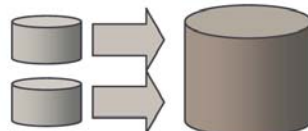
Performance Management Strategy

- Design strategic management model
- Metrics definition
- Performance Management Roadmap



Data warehouse development

- Architecture design
- ETL design
- Data Quality Assurance
- ETL Development



Reporting en analysis

- Reporting
- Dashboard development
- Data mining and analysis

