

# Databaser fra et logikkperspektiv

Evgenij Thorstensen

IFI, UiO

Høst 2013

# Outline

- 1 Relasjonsdatabaser
- 2 Fra relasjonsdatabase til modell
- 3 Spørninger som formler
- 4 Query containment

# Relasjonsdatabaser

Litt uformelt er en relasjonsdatabase en samling tabeller.

Brukernavn	Navn	FDato	Stillingskode
evgenit	Evgenij Thorstensen	05.07.1987	SKO1352

Kode	Tittel	Ansvarlig	Semester
INF3170	Logikk	arild	H

# Relasjonsdatabaser

Litt uformelt er en relasjonsdatabase en samling tabeller.

Brukernavn	Navn	FDato	Stillingskode
evgenit	Evgenij Thorstensen	05.07.1987	SKO1352

Kode	Tittel	Ansvarlig	Semester
INF3170	Logikk	arild	H

Informasjonen hentes via spørninger, for eksempel  
“SELECT Kode, Navn FROM Kurs, Personer WHERE  
Ansvarlig=Brukernavn AND Brukernavn='arild';”

Kode	Navn
INF3170	Arild Waaler

# Skjema

Mer presist, så består en relasjonsdatabase av to deler, et *skjema* og en *instans*.

Skjemaet er en samling relasjonssymboler med navngitte attributter, hver med sitt domene.

Personer(Brukernavn, Navn, FDato, Stillingskode)

Kurs(Kode, Tittel, Ansvarlig, Semester)

Domenet lar vi være implisitt; typisk int, string, etc.

## Instans

Gitt et skjema  $\Sigma$ , så er en instans over  $\Sigma$  en mengde tupler (rader) for hvert relasjonssymbol i skjemaaet, med riktig aritet og domener.

Vi kaller elementene i tuplene verdier.

Hvis vi vil, kan vi dropper navngitte attributter, og bruke indekser istedenfor (aritet må angis).

Dette likner jo veldig på noe dere har sett før.

## Fra database til modell

Vi kan tolke et skjema som en signatur, ved å la relasjonene være predikatsymboler.

Så kan vi tolke en instans som en modell  $\mathcal{M}$ , på følgende måte:

- Alle verdier er konstanter, og tolkes som seg selv.
- For predikatsymboler lar vi  
 $P^{\mathcal{M}} = \{\langle c_1^{\mathcal{M}}, \dots, c_n^{\mathcal{M}} \rangle \mid \langle c_1, \dots, c_n \rangle \text{ er i relasjonen til } P\}.$

La oss skrive  $DB^{\mathcal{M}}$  for modellen vi får fra databasen DB.

## Fra spørninger til formler

Siden en database svarer til en (endelig) modell, så bør formler svare til spørninger.

Gitt en lukket førsteordens formel  $\phi$  og en database DB, så kan vi sjekke om  $DB^{\mathcal{M}} \models \phi$ .

Dette er en *boolsk* spørring, “finnes det noe i DB som passer?”.

## Fra spørninger til formler

Siden en database svarer til en (endelig) modell, så bør formler svare til spørninger.

Gitt en lukket førsteordens formel  $\phi$  og en database DB, så kan vi sjekke om  $DB^M \models \phi$ .

Dette er en *boolsk* spørring, "finnes det noe i DB som passer?".

For eksempel, har alle kurs en ansvarlig som står i persontabellen?

$\forall$  Kode, Tittel, Ansvarlig, Semester

Kurs(Kode, Tittel, Ansvarlig, Semester)  $\rightarrow$

$\exists$  Navn, FDato, Stillingskode

Personer(Ansvarlig, Navn, FDato, Stillingskode)

$\forall x, y, z, v. Kurs(x, y, z, v) \rightarrow \exists x', y', z'. Personer(z, x', y', z')$

## Spørninger som returnerer data

SQL-spørringen i begynnelsen returnerer data. Slike spørninger svarer til formler med frie variable.

Vi husker at en variabel er *fri* i en formel hvis den ikke er innefor skopet til en kvantor. Vi skriver  $FV(\phi)$  for mengden av de frie variablene i  $\phi$ .

For eksempel, hvis

$$\phi = \forall x, y, v. Kurs(x, y, z, v) \rightarrow \exists y', z'. Personer(z, x', y', z'), \text{ så er}$$
$$FV(\phi) = \{z, x'\}$$

Intuitivt, finn alle  $z, x'$  slik at...

## Spørninger med frie variabler

La DB være en database over  $\Sigma$ , og  $\phi$  en formel med  $FV(\phi) = \{x_1, \dots, x_k\}$ .

Vi definerer svarene til  $\phi$  over DB som mengden av tupler  $\langle c_1, \dots, c_k \rangle$  med konstanter fra  $\Sigma$  slik at  $DB^{\mathcal{M}} \models \phi[c_1/x_1, \dots, c_k/x_k]$ .

Vi skriver  $Ans(\phi, DB)$  for denne mengden.

## Spørninger med frie variabler

La DB være en database over  $\Sigma$ , og  $\phi$  en formel med  $FV(\phi) = \{x_1, \dots, x_k\}$ .

Vi definerer svarene til  $\phi$  over DB som mengden av tupler  $\langle c_1, \dots, c_k \rangle$  med konstanter fra  $\Sigma$  slik at  $DB^{\mathcal{M}} \models \phi[c_1/x_1, \dots, c_k/x_k]$ .

Vi skriver  $Ans(\phi, DB)$  for denne mengden.

La DB være  $\{Pab, Pbc\}$  og  $\phi = \exists y(Pxy \wedge Pyz)$ . Da er  $Ans(\phi, DB) = \{\langle a, c \rangle\}$ .

## SPJ-spørringer

En veldig vanlig type spørring i SQL er select-project-join.

```
SELECT v1, v2... FROM T1, T2... WHERE  
Attr1 = Attr2, Attr3 = Attr4...
```

Vi kombinerer flere tabeller med likhet mellom attributter.

Denne typen spørring kan vi enkelt karakterisere ved hjelp av det logiske perspektivet.

## Konjunktive spørninger

En *konjunktiv spørring* (CQ) er bygd opp av konjunksjoner av atomer og eksistensielle kvantorer. Formelt:

- Alle atomære formler er med i CQ.
- Hvis  $\phi$  og  $\psi$  er i CQ, så er  $\phi \wedge \psi$  i CQ.
- Hvis  $\phi$  er i CQ, så er  $\exists x\phi$  i CQ.

Vi tillater altså ikke negasjon, disjunksjon, etc.

Siden samme variabel alltid tolkes likt, har vi likhet:  $Pxy \wedge Pzv \wedge y = z$  er det samme som  $Pxy \wedge Pyv$ .

## Oppsummering så langt

Vi kan tolke databaser som endelige modeller, ved å la hver tabell bli et predikat.

Spørninger blir da formler med frie variable.

En spesiell type spørninger er SPJ, som tilsvarer eksistenskvantifiserte konjunktive formler.

# Optimisering av spørninger

Hvis vi har en komplisert spørring, eller flere spørringer å kjøre, kan deler være overflødige (redundant).

Da kan spørringen optimiseres. For eksempel er  $P_{xx} \wedge \exists y P_{xy}$  ekvivalent til  $P_{xx}$ .

Vi ønsker å oppdage slike ting uten å evaluere spørringene.

# Query containment

## Definisjon (Query containment)

La  $\phi$  og  $\psi$  være to spørninger slik at  $FV(\phi) = FV(\psi)$ . Vi sier at  $\phi$  er inneholdt (contained) i  $\psi$  hvis  $Ans(\phi, DB) \subseteq Ans(\psi, DB)$  for alle DB.

Med andre ord, svarene til  $\phi$  er alltid en del av svarene til  $\psi$ .

Hvis  $\phi \subseteq \psi$  og  $\psi \subseteq \phi$ , så er spørringene ekvivalente.

Med et logisk perspektiv kan vi lett finne en måte å teste query containment på.

# Query containment, semantikk

La  $\phi$  og  $\psi$  være to lukkede formler.

## Theorem

*Vi har  $\phi \subseteq \psi$  hvis og bare hvis  $\phi \models \psi$ .*

Hvorfor det? Jo, fordi alle databaser = alle endelige modeller.

$\phi \subseteq \psi$  betyr at for alle  $DB^{\mathcal{M}}$ , hvis  $DB^{\mathcal{M}} \models \phi$ , så  $DB^{\mathcal{M}} \models \psi$ . Det er definisjonen av logisk konsekvens.

Hva hvis vi har frie variable?

## QC med frie variable

Hvis vi har frie variable, kan vi bruke *universell tillukning*.

La  $\phi$  være en formel med frie variable  $FV(\phi) = \{x_1, \dots, x_k\}$ . Den universelle tillukningen av  $\phi$  er  $\forall x_1, \dots, x_k \phi$ .

La  $\phi$  og  $\psi$  være to formler med  $FV(\phi) = FV(\psi) = \{x_1, \dots, x_k\}$ .

### Theorem

*Vi har  $\phi \subseteq \psi$  hvis og bare hvis  $\forall x_1, \dots, x_k \phi \models \forall x_1, \dots, x_k \psi$ .*

## QC med frie variable, bevis del en

La  $\phi$  og  $\psi$  være to formler med  $FV(\phi) = FV(\psi) = \{x_1, \dots, x_k\}$ .

### Theorem

*Vi har  $\phi \subseteq \psi$  hvis og bare hvis  $\forall x_1, \dots, x_k \phi \models \forall x_1, \dots, x_k \psi$ .*

Hvorfor det? Vel, bare hvis-delen funker greit. Anta  $\phi \subseteq \psi$ . Hvis  $DB^{\mathcal{M}} \models \forall x_1, \dots, x_k \phi$ , så er  $Ans(\phi, DB)$  alle mulige k-tupler.

Siden  $\phi \subseteq \psi$ , så er alle disse også svarene til  $\psi$ , og da har vi at  $DB^{\mathcal{M}} \models \forall x_1, \dots, x_k \psi$ .

## QC med frie variable, bevis del to

La  $\phi$  og  $\psi$  være to formler med  $FV(\phi) = FV(\psi) = \{x_1, \dots, x_k\}$ .

### Theorem

*Vi har  $\phi \subseteq \psi$  hvis og bare hvis  $\forall x_1, \dots, x_k \phi \models \forall x_1, \dots, x_k \psi$ .*

For hvis-delen, anta at det ikke stemmer. Da har vi at

$\forall x_1, \dots, x_k \phi \models \forall x_1, \dots, x_k \psi$ , men ikke at  $\phi \subseteq \psi$ . Da finnes en DB slik at  $Ans(\phi, DB) \not\subseteq Ans(\psi, DB)$ .

## QC med frie variable, bevis del to

La  $\phi$  og  $\psi$  være to formler med  $FV(\phi) = FV(\psi) = \{x_1, \dots, x_k\}$ .

### Theorem

*Vi har  $\phi \subseteq \psi$  hvis og bare hvis  $\forall x_1, \dots, x_k \phi \models \forall x_1, \dots, x_k \psi$ .*

For hvis-delen, anta at det ikke stemmer. Da har vi at

$\forall x_1, \dots, x_k \phi \models \forall x_1, \dots, x_k \psi$ , men ikke at  $\phi \subseteq \psi$ . Da finnes en DB slik at  $Ans(\phi, DB) \not\subseteq Ans(\psi, DB)$ .

Da kan vi ta  $\langle c_1, \dots, c_k \rangle \in Ans(\phi, DB) - Ans(\psi, DB)$ . Siden  $\langle c_1, \dots, c_k \rangle \in Ans(\phi, DB)$ , så har vi at  $DB^M \models \phi[c_1/x_1, \dots, c_k/x_k]$ , men at  $DB^M \not\models \psi[c_1/x_1, \dots, c_k/x_k]$ .

Ergo kan ikke  $\forall x_1, \dots, x_k \phi \models \forall x_1, \dots, x_k \psi$  stemme, og vi har en motsigelse.

## QC, algoritme

Teoremet gir oss følgende algoritme: For å sjekke  $\phi \subseteq \psi$ , sjekk  
 $\forall x_1, \dots, x_k \phi \vdash \forall x_1, \dots, x_k \psi$ .

## QC, algoritme

Teoremet gir oss følgende algoritme: For å sjekke  $\phi \subseteq \psi$ , sjekk  
 $\forall x_1, \dots, x_k \phi \vdash \forall x_1, \dots, x_k \psi$ .

### Theorem (Trakhtenbrot 1949)

*For førsteordens formler  $\phi$  og  $\psi$  er  $\phi \models \psi$  ikke avgjørbart, selv over endelige modeller.*

Hva med konjunktive spørninger?

# QC for konjunktive spørninger

## Theorem

For konjunktive spørninger  $\phi$  og  $\psi$  er  $\phi \sqsubseteq \psi$  avgjørbart, og NP-komplett.

For å vise at dette problemet er avgjørbart, skal vi bruke *homomorfier*.

Siden en konjunktiv spørring  $\phi$  bare inneholder eksistenskvantorer og konjunksjoner, kan vi skrive den som  $\exists x_1, \dots, x_n \phi$ , hvor  $\phi$  er åpen.

La  $Q = \exists x_1, \dots, x_n \phi$  og  $R = \exists y_1, \dots, y_m \psi$  være to konjunktive spørninger med  $FV(Q) = FV(R)$ . En *homomorfi* fra  $Q$  til  $R$  er en substitusjon  $\theta$  slik at  $\phi\theta$  er en delformel av  $\psi$ .

Homomorfien skal ikke endre de frie variablene i  $Q$  og  $R$ !

## Homomorfiteoremet del en

### Theorem (Homomorfiteoremet)

La  $Q = \exists y_1, \dots, y_n \phi$  og  $R = \exists y_1, \dots, y_m \psi$  være konjunktive spørninger med  $FV(Q) = FV(R) = \{x_1, \dots, x_k\}$ . Vi har  $Q \subseteq R$  hvis og bare hvis det finnes en homomorfi fra  $R$  til  $Q$ .

Hvis-delen: Anta at vi har en homomorfi  $\theta$ . La DB være en database, og  $\langle c_1, \dots, c_k \rangle \in Ans(Q, DB)$ . Da finnes det elementer  $d_1, \dots, d_n$  i  $DB^{\mathcal{M}}$  slik at  $DB^{\mathcal{M}} \models \phi[c_1/x_1, \dots, c_k/x_k, d_1/y_1, \dots, d_n/y_n]$ .

Siden  $\psi\theta$  er en delformel av  $\phi$ , og begge er konjunksjoner av atomer, så må  $DB^{\mathcal{M}} \models \psi\theta[c_1/x_1, \dots, c_k/x_k, d_1/y_1, \dots, d_n/y_n]$ .

Ergo er  $\langle c_1, \dots, c_k \rangle \in Ans(R, DB)$ .

## Homomorfiteoremet del to

### Theorem (Homomorfiteoremet)

La  $Q = \exists y_1, \dots, y_n \phi$  og  $R = \exists y_1, \dots, y_m \psi$  være konjunktive spørninger med  $FV(Q) = FV(R) = \{x_1, \dots, x_k\}$ . Vi har  $Q \subseteq R$  hvis og bare hvis det finnes en homomorfi fra  $R$  til  $Q$ .

Bare hvis-delen: La oss lage en database  $Q_{DB}$  ved å la hvert atom av  $\phi$  være et tuppel i sin tabell. For eksempel,  $Pxy \wedge Pxx$  blir tuplene  $\langle x, x \rangle$  og  $\langle x, y \rangle$  i  $P$ .

## Homomorfiteoremet del to

### Theorem (Homomorfiteoremet)

La  $Q = \exists y_1, \dots, y_n \phi$  og  $R = \exists y_1, \dots, y_m \psi$  være konjunktive spørninger med  $FV(Q) = FV(R) = \{x_1, \dots, x_k\}$ . Vi har  $Q \subseteq R$  hvis og bare hvis det finnes en homomorfi fra  $R$  til  $Q$ .

Bare hvis-delen: La oss lage en database  $Q_{DB}$  ved å la hvert atom av  $\phi$  være et tuppel i sin tabell. For eksempel,  $Pxy \wedge Pxx$  blir tuplene  $\langle x, x \rangle$  og  $\langle x, y \rangle$  i  $P$ .

Hva er  $Ans(Q, Q_{DB})$ ? Jo,  $\{\langle x_1, \dots, x_k \rangle\}$ . Siden  $Q \subseteq R$ , må  $\langle x_1, \dots, x_k \rangle \in Ans(R, Q_{DB})$ . Men da finnes det konstanter  $c_1, \dots, c_m$  i  $Q_{DB}$  slik at  $Q_{DB}^M \models \psi[c_1/y_1, \dots, c_m/y_m]$ .

Vi skal vise at  $[c_1/y_1, \dots, c_m/y_m]$  er en homomorfi fra  $R$  til  $Q$ .

## Homomorfiteoremet del to, fortsettelse

### Theorem (Homomorfiteoremet)

La  $Q = \exists y_1, \dots, y_n \phi$  og  $R = \exists y_1, \dots, y_m \psi$  være konjunktive spørninger med  $FV(Q) = FV(R) = \{x_1, \dots, x_k\}$ . Vi har  $Q \subseteq R$  hvis og bare hvis det finnes en homomorfi fra  $R$  til  $Q$ .

Hva er  $\text{Ans}(Q, Q_{DB})$ ? Jo,  $\{\langle x_1, \dots, x_k \rangle\}$ . Siden  $Q \subseteq R$ , må  $\langle x_1, \dots, x_k \rangle \in \text{Ans}(R, Q_{DB})$ . Men da finnes det konstanter  $c_1, \dots, c_m$  i  $Q_{DB}$  slik at  $Q_{DB}^M \models \psi[c_1/y_1, \dots, c_m/y_m]$ .

Disse konstantene er termer fra  $Q$ . Siden  $\psi$  er en konjunksjon av atomer, betyr dette at  $Q_{DB}$  inneholder hvert atom i  $\psi[c_1/y_1, \dots, c_m/y_m]$ .

Pga. måten vi konstruerte  $Q_{DB}$  på, betyr dette at også  $Q$  inneholder alle disse atomene, og ergo er  $\psi[c_1/y_1, \dots, c_m/y_m]$  en delformel av  $\phi$ .

## QC for konjunktive spørninger, kompleksitet

Siden det er endelig mange homomorfier, er problemet avgjørbart.

Siden det å sjekke om en substitusjon er en homomorfi kan gjøres i lineær tid, er vi i klassen NP.

NP-hardhet følger av hardhet for homomorfi-problemet.

## QC for konjunktive spørninger, konklusjon

For konjunktive spørninger  $\phi$  og  $\psi$  med like frie variable er følgende påstander ekvivalente:

- $\phi \subseteq \psi$
- $\phi \models \psi$
- Det finnes en homomorfi fra  $\psi$  til  $\phi$

Den i midten viste vi for arbitrære formler.

Den siste følger av homomorfiteoremet.

## Databaser som spørninger

Det å representere en konjunktiv spørring som en database kan snus.

Vi kan representere en database DB som en boolsk konjunktiv spørring,  $DB_Q$ , ved lage en stor konjunksjon av alle tupler i alle relasjoner.

$$DB_Q = \bigwedge_{P \text{ relasjon i DB}} \bigwedge_{\langle c_1, \dots, c_n \rangle \in P} P(c_1, \dots, c_n)$$

# Databaser som spørninger

Det å representere en konjunktiv spørring som en database kan snus.

Vi kan representere en database DB som en boolsk konjunktiv spørring,  $DB_Q$ , ved lage en stor konjunksjon av alle tupler i alle relasjoner.

$$DB_Q = \bigwedge_{P \text{ relasjon i } DB} \bigwedge_{\langle c_1, \dots, c_n \rangle \in P} P(c_1, \dots, c_n)$$

Da kan vi vise følgende: Hvis  $\phi$  er en boolsk konjunktiv spørring, så har vi  $DB^M \models \phi$  hvis og bare hvis  $DB_Q \subseteq \phi$ .

# Databaser som spørninger

Det å representere en konjunktiv spørring som en database kan snus.

Vi kan representere en database DB som en boolsk konjunktiv spørring,  $DB_Q$ , ved lage en stor konjunksjon av alle tupler i alle relasjoner.

$$DB_Q = \bigwedge_{P \text{ relasjon i } DB} \bigwedge_{\langle c_1, \dots, c_n \rangle \in P} P(c_1, \dots, c_n)$$

Da kan vi vise følgende: Hvis  $\phi$  er en boolsk konjunktiv spørring, så har vi  $DB^M \models \phi$  hvis og bare hvis  $DB_Q \subseteq \phi$ .

Bevis: Ukeopgave!

## Bonus: NP-hardhet for konjunktive spørninger

### Definisjon (Trefargeproblemet)

Gitt en urettet graf  $G$ , kan den farges med tre farger slik at nabonoder har ulik farge?

Vi skal lage en spørring og en database som uttrykker dette problemet.

## Bonus: NP-hardhet for konjunktive spøringer

### Definisjon (Trefargeproblemet)

Gitt en urettet graf  $G$ , kan den farges med tre farger slik at nabonoder har ulik farge?

Vi skal lage en spørring og en database som uttrykker dette problemet.

Databasen inneholder en binær relasjon  $E$  med seks tupler:  $\langle r, g \rangle$ ,  $\langle r, b \rangle$ , og  $\langle b, g \rangle$ , samt deres symmetrier.

Spørringen lager vi ut fra  $G$ . Hver node er en variabel, og hver kant  $\langle v, w \rangle$  er et atom  $E(v, w)$ . Formelen blir da

$$\exists V \bigwedge_{\langle v, w \rangle \in E} E(v, w)$$