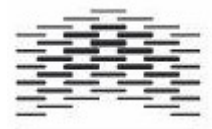


Hva er datakvalitet?

Hvordan skal arkivtjenesten forholde seg til det?

Thomas Sødring
Høyskolen i Oslo og Akershus
thomas.sodring@hioa.no
99570472



Bakgrunn

- Flere og flere depot institusjoner gjør seg klare til å ta imot elektroniske avleveringer
 - Hva vet depot om disse avleveringer?
 - Hvordan skal disse forvaltes? DIAS
 - Er det i det helt tatt mulig å bruke disse avleveringer?
- Flere stiller spørsmålet «*Kan vi måle kvaliteten på elektronisk materiale?*»

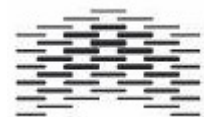
Datakvalitet

- Datakvalitet som et eget forskningsfelt har eksistert siden 1970 tallet
 - men det var etter 2000 tallet at flere og flere ble interessert i fagfeltet
 - Dette pga en eksplosjon i mengden av elektronisk data som ble generert
 - Hvordan dataene ble (og fortsatt blir) *håndtert på en ustrategisk* måte av mange selskaper

Hva er datakvalitet?

- Vi prøver på en definisjon

Datakvalitet angir i hvilken grad data i et system er i overensstemmelse med det virkelige scenarioet dataen representerer og er brukbar

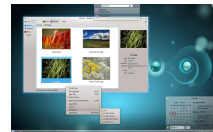
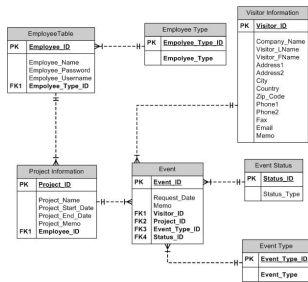


Systemutvikling

modellering

implementasjon

bruk



Kunder					
KunderNr	Fornavn	Etternavn	Adresse	Postnr	Sted
1	Mona	Lie	Storgata 4	0182	Oslo
2	Thomas	Hansen	Bakken 8b	1406	Skii
3	Eli	Rørvik	Saturningen 47	1808	Askim

Kunder					
KunderNr	Fornavn	Etternavn	Adresse	Postnr	Sted
1	Mona	Lie	Storgata 4	0182	Oslo
2	Thomas	Hansen	Bakken 8b	1406	Skii
3	Eli	Rørvik	Saturningen 47	1808	Askim

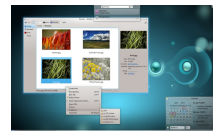
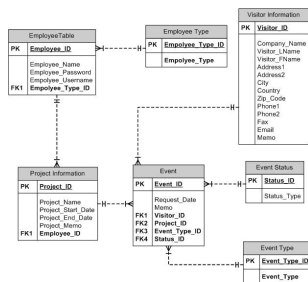
Hva er datakvalitet?

Datakvalitet angir i hvilken grad ***data i et system*** er i ***overensstemmelse*** med det virkelige ***scenariot dataen representerer*** og er ***brukbar***

modellering

implementasjon

bruk



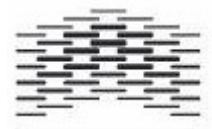
Kunder					
KunderNr	Fornavn	Efternavn	Adresse	Postnr	Sted
1	Mona	Lie	Storgata 4	0182	Oslo
2	Thomas	Hansen	Bakken 8b	1406	Skil
3	Eli	Rørvik	Saturnringen 47	1808	Askim

Kunder					
KunderNr	Fornavn	Efternavn	Adresse	Postnr	Sted
1	Mona	Lie	Storgata 4	0182	Oslo
2	Thomas	Hansen	Bakken 8b	1406	Skil
3	Eli	Rørvik	Saturnringen 47	1808	Askim



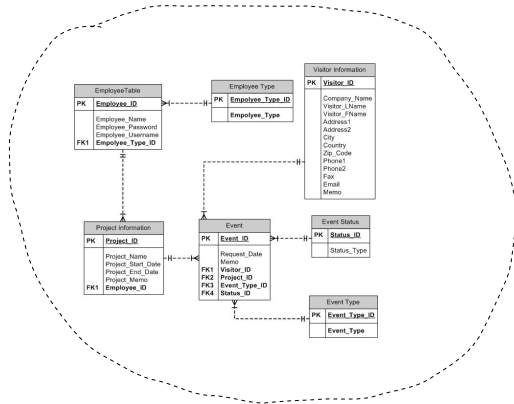
Hvordan oppstår dårlig datakvalitet

- Manglende datakvalitet er ofte et resultat av
 - *dårlig system design* eller
 - Scenarioet system representerer er ikke tilstrekkelig modellert
 - forbundet med *dårlige prosedyrer* ved data innførsel
 - Da er datakvalitet er en form av ***god arkivdanning!***



Hva er datakvalitet?

- *dårlig system design*
- *dårlige prosedyrer ved data innførsel*

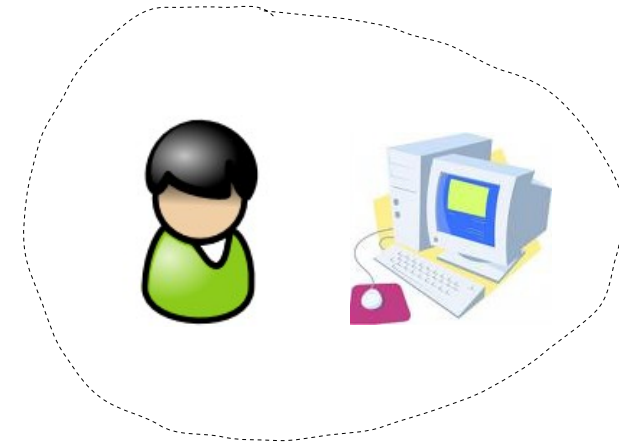


dårlig system design



Kunder					
KunderNr	Fornavn	Etternavn	Adresse	Postnr	Sted
1	Mona	Lie	Storgata 4	0182	Oslo
2	Thomas	Hansen	Bakken 8b	1406	Skil
3	Eli	Rørvik	Saturnringen 47	1808	Askim

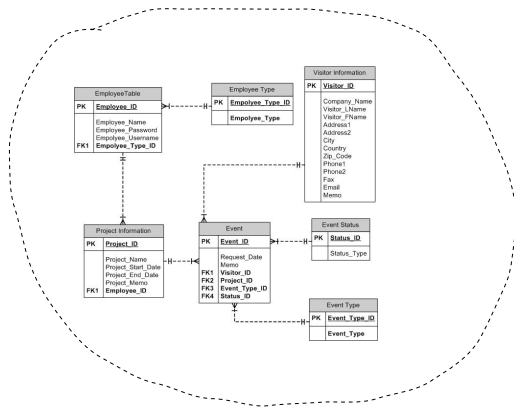
Kunder					
KunderNr	Fornavn	Etternavn	Adresse	Postnr	Sted
1	Mona	Lie	Storgata 4	0182	Oslo
2	Thomas	Hansen	Bakken 8b	1406	Skil
3	Eli	Rørvik	Saturnringen 47	1808	Askim



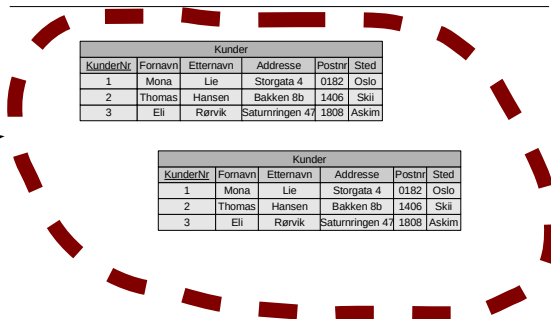
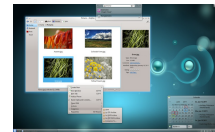
dårlige prosedyrer

Hva er datakvalitet?

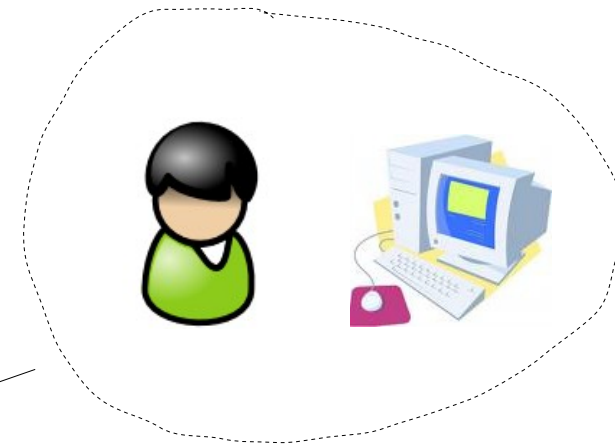
- *dårlig system design*
- *dårlige prosedyrer ved data innførsel*



dårlig system design



*Datakvalitetsproblemet
manifesterer seg her*



dårlige prosedyrer

*måles også her!
(til en hvis grad)*

DK ved danning eller bevaring?

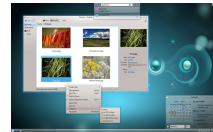
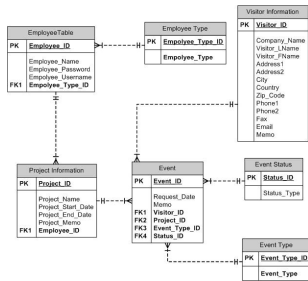
- Er datakvalitet noe som vi skal være opptatt av under danning eller bevaring?
- Hvordan oppleves datakvalitets problematikken i hver av fasene?
- Når depot overtar et uttrekk er det praktisk talt umulig å rette på kvaliteten
- Hvorfor tar ingen ansvar for datakvalitet?
 - Det står ikke noe om det i Noark standarden
 - Data er låst i leverandørenes systemer

DK ved danning eller bevaring?

modellering

implementasjon

bruk



Kunder					
KunderNr	Fornavn	Etternavn	Adresse	Postnr	Sted
1	Mona	Lie	Storgata 4	0182	Oslo
2	Thomas	Hansen	Bakken 8b	1406	Skil
3	Eli	Rørvik	Saturnringen 47	1808	Askim

Kunder					
KunderNr	Fornavn	Etternavn	Adresse	Postnr	Sted
1	Mona	Lie	Storgata 4	0182	Oslo
2	Thomas	Hansen	Bakken 8b	1406	Skil
3	Eli	Rørvik	Saturnringen 47	1808	Askim

Et uttrekk blir skapt fra data i databasen

DK ved danning eller bevaring?

- Depot vil arve DK fra dannelsesystemene
 - DK er noe depot bør være opptatt av
 - Danning har jo bare en 5 års perspektiv på kvalitet. Bevaring har ??
- Ja takk begge deler!
 - Mål og fiks datakvalitet underveis, la DK inngå i avleveringer
 - Kan depot nekte å ta imot et uttrekk pga dårligkvalitet?
 - Hvis den validerer mot DTD/XSD?

Nekte å ta i mot avleveringer?

- Statlig nivå
 - Riksarkivet vil sjekke at filene valideres mot en skjema
 - Teste med arkn4 / URD
 - Hvilken hjemmel kan Riksarkivet bruke til å nekte å ta imot en avlevering med dårlig datakvalitet
- IKA nivå
 - IKA er eid av de som skal levere, vanskeligere å nekte å ta imot
 - IKA bør måle kvalitet

Datakvalitet

- Datakvalitet måles i datakvalitetsdimensjoner
- Det finnes mange datakvalitetsdimensjoner. Ofte grupperer de:
 - Subjektive
 - Objektive
 - Prosess (ser ikke på det)
- De kan også grupperes på forskjellige måter*

*2 artikler

<http://www.emeraldinsight.com/journals.htm?articleid=841292&show=html>

<http://citeseer.ist.psu.edu/viewdoc/summary?doi=10.1.1.144.4212>

Subjektivt Datakvalitet

- Subjektive datakvalitets dimensjoner ofte basert på intervjueteknikker og svarene lar seg ofte ikke generalisere
- En saksbehandler kan ha en forståelse at dataen i saksbehandlingssystemet er av høy kvalitet utifra sitt perspektiv som bygger på mange års erfaring og bruk av systemet
 - En annen saksbehandler kan se på de samme dataene og mene at kvaliteten er dårlig
- Et eksempel av denne typen datakvalitet er forståelighets dimensjonen, dvs hvor enkelt det er å forstå informasjonen

Objektivt Datakvalitet

- Objektive datakvalitets dimensjoner vurderes gjennom en uavhengig analyse av dataen i systemet (databasen) og er ofte basert på programvare som går inn i databasen og måler kvaliteten i tabellene
 - Et eksempel på en slik dimensjon er fullstendighet eller om all relevant data er innhentet og lagret
 - Med uavhengighet menes hvis målingen kjøres to ganger over samme database vil du få samme resultat

Er data mellom systemer konsistent?
Eksister det duplikater?

konsistens

Er all data som trengs tatt med?

fullstendighet

Gjenspeiler data virkeligheten?

korrekthet

Datakvalitet

integritet

Er referanser mellom entiteter og attributter konsistent?

tidsriktighet

Er data tilgjengelig når de skal være tilgjengelig?
Er data utdatert?

gyldighet

Kommer alle verdier fra domenen av verdier?

Datakvalitets dimensjoner

Prototype datakvalitets GUI

The screenshot shows a web browser window with the address bar displaying `http://ark1.hio.no:8180/n5.ui.dq/`. The browser tabs include 'Quality Indic...', 'JODConvert...', 'PyODConver...', 'anytopdf - C...', 'Batch comm...', 'Database m...', and 'CPSC 343:'. The browser's address bar shows the URL `http://ark1.hio.no:8180/n5.ui.dq/`. The browser's search bar contains 'Google Translate', 'Hello, World! Web Ap...', 'Document-3297.phtm...', 'EJB Workshop', 'Loading...', and 'JBoss web servic'. The main content area is titled 'Data Quality Monitor' and contains a table and a gauge chart.

Dimension	Case	Author
Inclusion Dependency	Case 1	J. Smith
Processing Delay	Case 2	James Jones
Disposal	Case 3	J. Jones
Completeness	Case 4	J. Brown
Syntactic Accuracy	Case 5	Sarah Jones
Consistency	Case 6	J. Smith
Functional Dependency	Case 7	J. Brown
Correctness	Case 8	James Smith
	Case 9	J. Jones
	Case 10	Jack Jones

The gauge chart is a circular meter with a scale from 0.0 to 100.0. The needle is positioned at 0.0, and the text '0.0' is displayed in a white box at the bottom center of the gauge.

Datakvalitet og Noark

- Noark indirekte sørger for et minimum av datakvalitet
 - Det er ingen eksplosjon av elektronisk informasjon brukt på en ustrategisk måte i offentlig sektor
 - Eller?

Prosjekt om datakvalitet

- Kan vi *kvantifisere* «kvalitet» på en *objektiv måte*
 - Kan arkivfeltet kan dra nytte av datakvalitet
 - Subjektiv kan også være nyttig