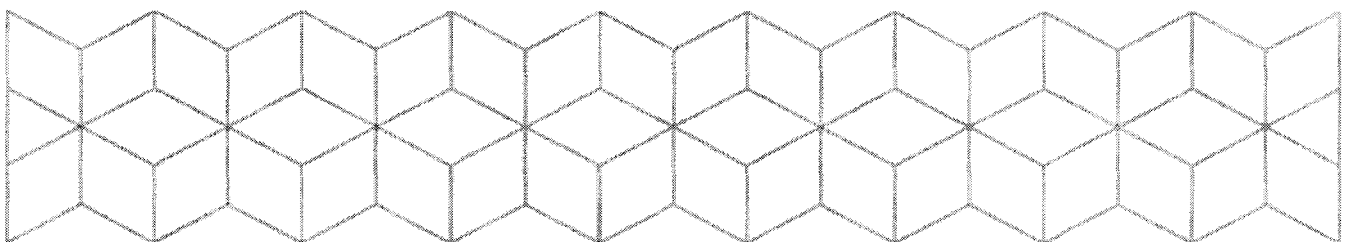


# EKSAMEN

<b>Emnekode:</b> ITF301416	<b>Emnenavn:</b> Store datamengder – analyse og prosessering
<b>Dato:</b> 16.12.16	<b>Eksamenstid:</b> 09:00 – 13:00
<b>Hjelpemidler:</b> Ingen	<b>Faglærer:</b> Edgar Bostrøm
<b>Om eksamensoppgaven og poengberegning:</b> <p>Oppgavesettet består av 4 sider inklusiv denne forsiden og et vedlegg.</p> <p>Kontroller at oppgaven er komplett før du begynner å besvare spørsmålene.</p> <p>Alle de 4 oppgavene skal besvares og teller likt ved sensurering.</p> <p>Les gjennom hele oppgavesettet med en gang, slik at du kan stille evt. spørsmål til hele settet når faglærer kommer.</p>	
<b>Sensurfrist:</b> 16.1.17 <p>Karakterene er tilgjengelige for studenter på Studentweb senest 2 virkedager etter oppgitt sensurfrist. <a href="http://www.hiof.no/studentweb">www.hiof.no/studentweb</a></p>	



# Oppgave 1. Relasjonsalgebra og optimalisering. Tid: 1 time.

Vi tar utgangspunkt i samme system som i et av prosjektene som ble gitt i høst.

## Prosjekt

Prosjektnr Prosjektnavn

## Ansatt

Ansattnr Fødselsnr Etternavn Fornavn

## ProsjektArbeide

Prosjektnr Ansattnr Timetall

- a) Skriv ut hvor mange timer den ansatte med fødselsnr 15129911111 har brukt på prosjektet «Renovering av Svinesundbrua». Følgende forslag virker. Forklar hva som gjøres, hvorfor det ikke er lurt å bruke dette relasjonsalgebrautsagnet.

$\pi_{\text{Timetall}}(\sigma_{\text{Fødselsnr} = 15129911111}(\sigma_{\text{Prosjektnavn} = \text{«Renovering av Svinesundsbrua»}))$

$\sigma_{\text{Prosjekt.Projektnr} = \text{ProsjektArbeide.Projektnr and Ansatt.Ansattnr} = \text{ProsjektAnsatt.Projektnr} (\text{Prosjekt X ProsjektArbeide X Ansatt}})$

- b) Lag en mest mulig effektiv formulering av samme spørsmål som over.

- c) Skriv ut projektnr og -navn på de prosjektene som ikke er satt i gang (dvs. prosjekter hvor det ikke finnes prosjektarbeide). Bruk helst mengdedifferanse, helst også semijoin.

- d) Skriv ut Ansattnr, Etternavn, Fornavn på alle ansatte, sammen med eventuelle prosjekter de har vært med på, med Projektnr og Timetall, slik som vist under. Tips: outer join. Vi antar at det finnes en **NULL** som kan brukes på samme måte som i SQL.

Ansattnr	Etternavn	Fornavn	Prosjektnr	Timetall
712	Andersen	Anne	P_171	40
712	Andersen	Anne	P_131	70
712	Andersen	Anne	P_151	30
399	Bindersen	Binder	NULL	NULL
821	Caspersen	Casper	P_131	98
540	Dauidsen	David	NULL	NULL

- e) Skriv ut Ansattnr, Etternavn og Fornavn på de ansatte som har deltatt i alle prosjekter som er satt i gang.

- f) I lærebøker og på nettet finner vi mange råd om «avoid using DISTINCT in SQL». Hvorfor? Og: er det en aktuell problemstilling i relasjonsalgebra? Hvorfor/hvorfor ikke?

## Oppgave 2. Replikasjon.

Tid: 1 time.

- a) Beskriv hva distribusjon/replikasjon av data er, og hvilke fordeler og ulemper det innebærer. (Tid: 30 min.)
- b) Forklar hva 2-fase-gjennomføring (2PC, 2 Two phase commit) av distribuerte transaksjoner er, og fordeler/begrensninger/alternativer til denne protokollen. (Tid: 15 min.)
- c) Det finnes et kjent system innenfor Big Data som bruker datadistribusjon. Hvordan virker dette? (Tid: 15 min.)

## Oppgave 3.

Tid: 1 time.

- a) Gi to gode grunner til fremveksten av Big Data teknologi
- b) Graf-databaser er veldig gode til en spesiell bruk, hva slags? Gi gjerne et eksempel på bruk av graf-databaseteknologi.
- c) Map Reduce er en sentral del av Hadoop. Gi et eksempel på problemer som Map Reduce er godt egnet for å løse, men også et eksempel på problemer Map Reduce ikke er egnet til å løse.
- d) Hvorfor er Apache Spark i utgangspunktet raskere enn Apache Hadoop.
- e) Sanntidsanalyse benytter seg typisk av tre tidsvinduer for å analysere data. Hva heter de og hvordan fungerer de?

## Oppgave 4.

Tid: 1 time.

- a) Beskriv hovedelementene i en typisk Big Data Analysis / Data Science prosess.
- b) Hvorfor kan det være viktig å rense data for outliers?
- c) Beskriv forskjellen mellom Bias-feil og Variance-feil.
- d) Tenk deg at vi har et klassifiseringsproblem med to klasser, A og B, og vi har 99000 eksempler av klasse A og 1000 eksempler av klasse B i vårt datasett. Beskriv hvordan du vil gå frem for å lage og evaluere en klassifiseringsmodell, og hvordan du vil gå frem for å sammenligne to alternative modeller.
- e) Beskriv hvordan k-nearest-neighbours modellen fungerer.

## VEDLEGG: Relasjonsalgebra - vanlige operasjoner.

<b>Mengdeoperasjoner:</b>	<i>Notasjon, variant 1</i>	<i>Notasjon, variant -2</i>
Union	$R \cup S$	R union S
Snitt	$R \cap S$	R intersect S
Mengdedifferanse	$R - S$ $R \setminus S$	R difference S R minus S
Mengdeprodukt, kartesisk produkt ("alle mot alle")	$R \times S$	R product S R times S
<b>Spesielt for relasjoner:</b>		
Horisontalt utvalg (sigma)	$\sigma_{\langle \text{betingelse} \rangle}(R)$	R where <bet.> R where <bet.>
Vertikalt utvalg (pi)	$\pi_{\langle \text{attributtliste} \rangle}(R)$	R[<attributtliste>]
Mengdedivisjon. (Gitt R[c,d] og S[d]. c er med i mengden R dividert med S hvis c i R forekommer sammen med alle d-er som finnes i S.)	$R \div S$ $R / S$	R divideby S
<b>Spesialiteter av produkt:</b>		
$\theta$ -join (produkt med en eller annen betingelse på kompatible attributter, f.eks. >, <, og komb.)	$R \bowtie_{\langle \text{betingelse} \rangle} S$	R join<betingselse> S (R join S) where <bet.>
Equi-join ( $\theta$ -opersjonen er =)	" "	" "
Natural join (Equi-join hvor felles attributt kommer bare en gang) ** den mest vanlige jointypen **	" "	" "
<b>Varianter for produkt:</b>		
Outer join, normalt venstre.. (alle i R, samt alle fra S som oppfyller koblingsbetingelsen)	$R \bowtie_{\langle \text{bet.} \rangle} S$	R left join<bet.> S
Full join (alle i R, alle i S, samt alle som oppfyller koblingsbet.)	$R \bowtie_{\langle \text{bet.} \rangle} S$	R full join<bet.> S
Semijoin (de i R som tilfredsstillers R join<betingselse> S)	$R \triangleright_{\langle \text{betingelse} \rangle} S$	R semijoin<bet.> S

Dersom betingelsen er på primær/fremmednøkkelkombinasjoner, droppes ofte <betingelse>.